

500,112

(19) 世界知的所有権機関
国際事務局



(43) 国際公開日
2003 年 7 月 17 日 (17.07.2003)

PCT

(10) 国際公開番号
WO 03/058500 A1

(51) 国際特許分類: G06F 17/30,
C12Q 1/68, C12N 15/00, G01N 33/50

(21) 国際出願番号: PCT/JP03/00011

(22) 国際出願日: 2003 年 1 月 6 日 (06.01.2003)

(25) 国際出願の言語: 日本語

(26) 国際公開の言語: 日本語

(30) 優先権データ:
特願 2001-402081

2001 年 12 月 28 日 (28.12.2001) JP

(71) 出願人 (米国を除く全ての指定国について): セレス
ター・レキシコ・サイエンス株式会社 (CELESTAR
LEXICO-SCIENCES, INC.) [JP/JP]; 〒261-8501 千葉
県 千葉市美浜区 中瀬 1 丁目 3 番地 幕張テクノガ
ーデン D 1 7 Chiba (JP).

(72) 発明者; および

(75) 発明者/出願人 (米国についてのみ): 上村 泰央 (UE-
MURA, Yasuo) [JP/JP]; 〒261-8501 千葉県 千葉市美浜

区 中瀬 1 丁目 3 番地 幕張テクノガーデン D 1 7 セ
レスター・レキシコ・サイエンス株式会社内 Chiba
(JP). 蓬萊 尚幸 (HORAI, Hisayuki) [JP/JP]; 〒261-8501
千葉県 千葉市美浜区 中瀬 1 丁目 3 番地 幕張テクノ
ガーデン D 1 7 セレスター・レキシコ・サイエン
ス株式会社内 Chiba (JP).

(74) 代理人: 酒井 宏明, 外 (SAKAI, Hiroaki et al.); 〒
100-0013 東京都 千代田区 霞が関三丁目 2 番 6 号 東
京倶楽部ビルディング Tokyo (JP).

(81) 指定国 (国内): US.

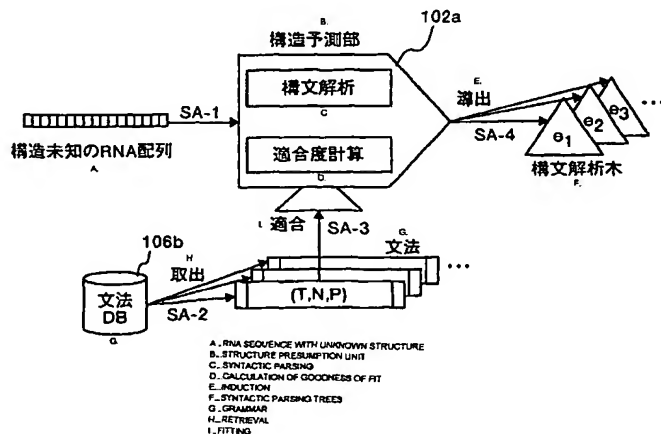
(84) 指定国 (広域): ヨーロッパ特許 (AT, BE, BG, CH, CY,
CZ, DE, DK, EE, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL,
PT, SE, SK, TR).

添付公開書類:
— 国際調査報告書

2 文字コード及び他の略語については、定期発行される
各 PCT ガゼットの巻頭に掲載されている「コードと略語
のガイダンスノート」を参照。

(54) Title: RNA SEQUENCE ANALYZER, AND RNA SEQUENCE ANALYSIS METHOD, PROGRAM AND RECORDING MEDIUM

(54) 発明の名称: RNA 配列解析装置、RNA 配列解析方法、プログラム、および、記録媒体



(57) Abstract: An RNA sequence analyzer which is provided with a grammar storage unit wherein the structural topology of an RNA secondary structure and the corresponding generative grammar fitting for the topology are stored; a syntactic parsing unit wherein an RNA sequence is applied to the generative grammar to induce syntactic parsing trees; and a unit for calculating goodness of fit wherein the goodness of fit for the syntactic parsing trees induced by the syntactic parsing unit is calculated.



(57) 要約:

本発明にかかるRNA配列解析装置は、RNA二次構造の構造トポロジーと当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段と、上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段とを備えている。

明 細 書

RNA配列解析装置、RNA配列解析方法、プログラム、および、記録媒体

5 技術分野

本発明は、RNA配列解析装置、RNA配列解析方法、プログラム、および、記録媒体に関し、特に、RNAの二次構造を予測し、また、DNA配列から遺伝子部分を予測するRNA配列解析装置、RNA配列解析方法、プログラム、および、記録媒体に関する。

10

背景技術

RNA配列は、A（アデニン）、C（シトシン）、G（グアニン）、U（ウラシル）の4種の塩基により構成されるが、その一部は回文言語のような入れ子の状態となり、相補的な塩基同士（AとU、GとC、稀にGとU）が結合して二次構造を構成している。RNA配列の二次構造は、第1図に示すように、多種の構造トポロジーを有している。ここで、積み重ねられた塩基対の連続領域をステム（stem）と呼び、また、塩基対で挟まれた一本鎖の部分配列をループ（loop）と呼ぶ。ステムの端のループをヘアピンループという（第1図のa）。ステムの中にある一本鎖の塩基は、それがステムの片側だけにあるときバルジループ（bulge loop）と呼び（第1図のb）、ステムの両側にあるときは内側ループ（internal loop）と呼ぶ（第1図のc）。放射状に3個以上のステムが出ているものをマルチループ（multi-branched loop）と呼ぶ。また、入れ子ではない塩基対があるとき、シュードノット（pseudoknots）と呼ぶ（第1図のd）。

25 ここで、RNA配列を形式文法（生成文法）により構文解析することによりその二次構造を予測する手法が存在するが、正規文法では回文言語を記述することができないため、一般に、RNA二次構造解析においては、木文法（tree

adjoining grammars) や、文脈自由文法 (CFG) などを用いて構文解析を行い構造モデリング (構造トポロジー表現) を求める手法が考案されている。

例えば、Yasuo Uemura 等著「Tree adjoining grammars for RNA structure prediction (Theoretical Computer Science 210 1999 277p-303p)」 (以下「文献1」という) には、木文法による構造モデリングと、パーザ (parser) を利用したエネルギー極小化による RNA 二次構造予測方法が開示されている。

また、Elena Rivas and Sean R. Rddy 著「The language of RNA: a formal grammar that includes pseudoknots (BIOINFORMATICS vol. 16 no. 4 2000 334p-340p)」 (以下「文献2」という) には、Crossed-interaction Grammars などの独自の拡張を施した文脈自由文法 (CFG) による構造モデリングと、パーザを利用したエネルギー極小化による RNA 二次構造予測方法が開示されている。

また、Michael Zuker 著「Prediction of RNA Secondary Structure by Energy minimization (July 8, 1996)」 (以下「文献3」という) には、形式文法やパーザを用いず動的計画法 (Dynamic Programming) によって RNA 二次構造を予測する方法を用いた RNA 配列解析システムである Mfold (製品名) が開示されている。これらの文献では、形式文法や動的計画法などの手法と、エネルギー極小化手法とを組み合わせることによって RNA 二次構造予測精度を高めている。

第2図は、従来技術による RNA の二次構造がステムループをとる場合の構文解析木の一例を示す図である。第2図の a に示す RNA 配列の二次構造を第2図

のbに、また、構文解析木を第2図のcに示している。ここで、部分木 (s u b
t r e e) は、内部の節を根とする構文解析木の断片である。RNA二次構造の
構造トポロジーについて構文解析木を作成して構文解析を行うことにより二次構
造解析を行う技術が研究されており、主要な構造トポロジーに対する文法が既知
5 となっている。

第3図は、従来技術によるRNA二次構造の構造トポロジーについて、文法が
固定されるとそれに対応した構造トポロジーが規定される (逆もまた可) ことを
示す概念図である。ここで、生成文法 (以下単に「文法」という) は、有限個の
記号と、有限個の生成規則 P (p r o d u c t i o n r u l e) から成る。記
10 号には、抽象的な非終端記号 N (n o n t e r m i n a l s y m b o l) と、
観測文字列に実際に現れる終端記号 T (t e r m i n a l s y m b o l) の2
種類がある。終端記号 T は、RNA配列の場合にはA、T、G、Cの4文字であ
る。第3図に示すように、各構造トポロジーについてそれぞれ対応する文法を定
義することができる。

また、第4図は、従来技術である木文法パーザを用いて、既知の文法からRN
A配列の構文解析木を導出する場合の一例を示す図である。まず、構造未知のR
NA配列を木文法パーザに入力する。ここで、木文法パーザは、入力された既知
の木文法に従ってRNA配列の構文解析を行い構文解析木を導出する機能、およ
び、導出された構文解析木について、ループや、塩基対とその他の二次構造要素
20 のそれぞれの自由エネルギーの合計などを計算することにより平衡自由エネルギ
ー (ΔG) などの値を求める機能などを有する (文献1から3参照) 。

ここで、木文法パーザは必ずしも構文解析木を導出するわけではなく、入力し
たRNA配列が当該文法に適合しない場合 (パーズが成功しない場合) には構文
解析木を導出しない (すなわち、構文解析木は0個になる) 。また、複数の構
25 文解析木が導出された場合には、エネルギー計算の結果、極小の自由エネルギー
となる1つの構文解析木を選択する。このとき、木文法パーザは導出過程の各段
階において自由エネルギー極小な部分構造を見つけていくことができる。また、

本文法パーザはエネルギー準最適な構文も出力することができる。このように、本文法パーザは、構文解析（パース）の途中でエネルギー計算を行うことにより、高速化と精度向上を実現している。

5 しかしながら、従来の本文法パーザなどにより構文解析とエネルギー計算を行う手法を用いたRNA二次構造予測システムにおいては、RNA配列や抽出した文法を統合的に管理し、集積した文法やRNA配列を用いて二次構造予測などをより効率的に行うものは存在しなかったという問題点がある。

また、与えられた特定の二次構造を取り得るようなRNA配列を検索するような方法は存在しなかったという問題点がある。

10 また、複数のRNA配列に共通な二次構造を簡易に抽出する方法は存在しなかったという問題点がある。

また、RNA配列からRNA二次構造に基づく類似度を簡易に求める方法は存在しなかったという問題点がある。

15 さらに、DNA配列から遺伝子部分を発見するための手法としては、ホモロジー検索やモチーフ検索などを利用する手法が一般的であるが、未知の遺伝子部分の発見には利用できないという問題点がある。ここで、従来技術で説明したように、RNA配列の構造トポロジーを予測可能な生成文法が求められているが、既知の生成文法により導出された構文解析木を用いた遺伝子発見方法は存在しなかったという問題点がある。

20 このように、従来のシステム等は数々の問題点を有しており、その結果、システムの利用者および管理者のいずれにとっても、利便性が悪く、また、利用効率が悪くものであった。

従って、本発明は、RNA配列や抽出した文法を統合的に管理し、集積した文法やRNA配列を用いて二次構造予測や新たな解析手法などをより効率的に行う
25 ことのできる、RNA配列解析装置、RNA配列解析方法、プログラム、および、記録媒体を提供することを目的としている。

発明の開示

本発明にかかるRNA配列解析装置は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段と、上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、上記適合度計算手段により計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を上記適合度が高い順にソートするソート手段と、上記ソート手段によりソートされた上記構文解析木を上記RNA配列の二次構造の候補として出力する出力手段とを備えたことを特徴とする。

この装置によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算し、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を適合度が高い順にソートし、ソートされた構文解析木をRNA配列の二次構造の候補として出力するので、一配列に対して多文法で構文解析を行うことができるようになる。すなわち、各生成文法に対してそれぞれ構文解析し適合度計算を行い適合度を得る。その結果、生成文法ごとに適合度が得られることになり、それらの適合度をソートすることによって生成文法に順位を付ける。これにより、生成文法に対する構造トポロジーにも順位が付けられることになるので、最終的にRNA配列が取得する可能性の高い順に構造トポロジーを確認することができるようになる。

つぎの発明にかかるRNA配列解析装置は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段と、上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、上記適合度計算手段により計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力する出力手段

とを備えたことを特徴とする。

この装置によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力するので、多配列に対して一文法で構文解析を行うことができるようになる。すなわち、与えられた特定の構造トポロジーに対し、対応する生成文法を取得し、これを用いてRNA配列データベースに格納されているすべてまたは一部のRNA配列をそれぞれ構文解析し、ある閾値以下の適合度で構文解析に成功したRNA配列群を結果として出力する。これにより、与えられた特定の二次構造を取り得るようなRNA配列を検索することができるようになる。

つぎの発明にかかるRNA配列解析装置は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段と、上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、上記適合度計算手段により計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出手段と、上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出手段にて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分をマークすることにより、上記RNA配列間で共通に有する構造トポロジーを可視化する共通構造マトリックス作成手段とを備えたことを特徴とする。

この装置によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したR

NA配列を抽出し、構造トポロジーとRNA配列とを二次元マトリックスで表示し、二次元マトリックスにおいて抽出されたRNA配列と構造トポロジーに対応する格子部分をマークすることにより、RNA配列間で共通に有する構造トポロジーを可視化するので、RNA配列間の共通構造を容易に発見することができるようになる。

5 つぎの発明にかかるRNA配列解析装置は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、利用者が入力したDNA配列から転写されるRNA配列を作成するRNA配列作成手段と、上記RNA配列作成手段により作成された上記RNA配列に対して上記生成文法を適用して構文解析木を導出する構文解析手段と、上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、上記適合度計算手段により計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列に対応する上記DNA配列部分を遺伝子の候補として予測する遺伝子予測手段とを備えたことを特徴とする。

10 この装置によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、利用者が入力したDNA配列から転写されるRNA配列を作成し、作成されたRNA配列に対して生成文法を適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列に対応するDNA配列部分を遺伝子の候補として予測するので、DNA配列のうち既知のトポロジーを有する可能性のあるRNA配列に対応する部分について、遺伝子部分である可能性があることを予測することができるようになる。

25 つぎの発明にかかるRNA配列解析装置は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段

と、上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、上記適合度計算手段により計算された上記適合度に基づいて上記RNA配列間の類似度を計算する類似度計算手段とを備えたことを特徴とする。

5 この装置によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度に基づいてRNA配列間の類似度を計算するので、RNA構造の類似度を容易に求めることができるようになる。

10 つぎの発明にかかるRNA配列解析装置は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段と、上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、上記適合度計算手段により計算された上記適合度のうち予
15 め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出手段と、上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出手段にて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分に上記適合度を表示する適合度マトリックスを作成する適合度マトリックス作成手段と、上記適合度マト
20 リックス作成手段にて作成された上記適合度マトリックスについて、上記適合度により上記構造トポロジーをソートし、他のRNA配列について当該ソートされた上記構造トポロジーの順番に対応する上記生成文法により構文解析を行い上記適合度が最大となる上記構文解析木を求め、予め定めた条件を満たす上記適合度を持つ上記構文解析木に対応する上記他のRNA配列を抽出する共通構造抽出手
25 段とを備えたことを特徴とする。

 この装置によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構

文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を抽出し、構造トポロジーとRNA配列とを二次元マトリックスで表示し、二次元マトリックスにおいて抽出されたRNA配列と構造トポロジーに対応する格子部分に適合度を表示する適合度マトリックスを作成し、適合度マトリックスについて、適合度により構造トポロジーをソートし、他のRNA配列について当該ソートされた構造トポロジーの順番に対応する生成文法により構文解析を行い適合度が最大となる構文解析木を求め、予め定めた条件を満たす適合度を持つ構文解析木に対応する他のRNA配列を抽出するので、共通の構造を持つRNA配列を容易に発見することができるようになる。

また、本発明はRNA配列解析方法に関するものであり、本発明にかかるRNA配列解析方法は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を上記適合度が高い順にソートするソートステップと、上記ソートステップによりソートされた上記構文解析木を上記RNA配列の二次構造の候補として出力する出力ステップとを含むことを特徴とする。

この方法によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度を計算し、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を適合度が高い順にソートし、ソートされた構文解析木をRNA配列の二次構造の候補として出力するので、一配列に対して多文法で構文解析を行うことができるようになる。すなわち、各生成文法に対してそれぞれ構文解析し適合度計算を行い適合度を得る。

その結果、生成文法ごとに適合度が得られることになり、それらの適合度をソートすることによって生成文法に順位を付ける。これにより、生成文法に対する構造トポロジーにも順位が付けられることになるので、最終的にRNA配列が取り得る可能性の高い順に構造トポロジーを確認することができるようになる。

- 5 つぎの発明にかかるRNA配列解析方法は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力する出力ステップとを含むことを特徴とする。
- 10

- この方法によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力するので、多配列に対して一文法で構文解析を行うことができるようになる。すなわち、与えられた特定の構造トポロジーに対し、対応する生成文法を取得し、これを
- 15
- 20
- これを用いてRNA配列データベースに格納されているすべてまたは一部のRNA配列をそれぞれ構文解析し、ある閾値以下の適合度で構文解析に成功したRNA配列群を結果として出力する。これにより、与えられた特定の二次構造を取り得るようなRNA配列を検索することができるようになる。

- つぎの発明にかかるRNA配列解析方法は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合
- 25

度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出ステップと、上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出ステップにて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分をマークすることにより、上記RNA配列間で共通に有する構造トポロジーを可視化する共通構造マトリックス作成ステップとを含むことを特徴とする。

この方法によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を抽出し、構造トポロジーとRNA配列とを二次元マトリックスで表示し、二次元マトリックスにおいて抽出されたRNA配列と構造トポロジーに対応する格子部分をマークすることにより、RNA配列間で共通に有する構造トポロジーを可視化するので、RNA配列間の共通構造を容易に発見することができるようになる。

つぎの発明にかかるRNA配列解析方法は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、利用者が入力したDNA配列から転写されるRNA配列を作成するRNA配列作成ステップと、上記RNA配列作成ステップにより作成された上記RNA配列に対して上記生成文法を適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列に対応する上記DNA配列部分を遺伝子の候補として予測する遺伝子予測ステップとを含むことを特徴とする。

この方法によれば、RNA二次構造の構造トポロジーと、当該構造トポロジー

に適合する生成文法とを対応付けて格納し、利用者が入力したDNA配列から転写されるRNA配列を作成し、作成されたRNA配列に対して生成文法を適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列に対応するDNA配列部分を遺伝子の候補として予測するので、DNA配列のうち既知のトポロジーを有する可能性のあるRNA配列に対応する部分について、遺伝子部分である可能性があることを予測することができるようになる。

つぎの発明にかかるRNA配列解析方法は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度に基づいて上記RNA配列間の類似度を計算する類似度計算ステップとを含むことを特徴とする。

この方法によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度に基づいてRNA配列間の類似度を計算するので、RNA構造の類似度を容易に求めることができるようになる。

つぎの発明にかかるRNA配列解析方法は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出ステップと、上記構造トポロジーと上記RN

- A配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出ステップにて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分に上記適合度を表示する適合度マトリックスを作成する適合度マトリックス作成ステップと、上記適合度マトリックス作成ステップにて作成された上記適合度マトリックスについて、上記適合度により上記構造トポロジーをソートし、他のRNA配列について当該ソートされた上記構造トポロジーの順番に対応する上記生成文法により構文解析を行い上記適合度が最大となる上記構文解析木を求め、予め定めた条件を満たす上記適合度を持つ上記構文解析木に対応する上記他のRNA配列を抽出する共通構造抽出ステップとを含むことを特徴とする。
- 5
- 10 この方法によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を抽出し、構造トポロジーとRNA配列とを二次元マトリックスで表示
- 15 し、二次元マトリックスにおいて抽出されたRNA配列と構造トポロジーに対応する格子部分に適合度を表示する適合度マトリックスを作成し、適合度マトリックスについて、適合度により構造トポロジーをソートし、他のRNA配列について当該ソートされた構造トポロジーの順番に対応する生成文法により構文解析を行い適合度が最大となる構文解析木を求め、予め定めた条件を満たす適合度を持つ
- 20 つ構文解析木に対応する他のRNA配列を抽出するので、共通の構造を持つRNA配列を容易に発見することができるようになる。

また、本発明はRNA配列解析方法をコンピュータに実行させるプログラムに関するものであり、本発明にかかるプログラムは、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格

25 納ステップと、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算

された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を上記適合度が高い順にソートするソートステップと、上記ソートステップによりソートされた上記構文解析木を上記RNA配列の二次構造の候補として出力する出力ステップとを含むことを特徴とする。

- 5 このプログラムによれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算し、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を適合度が高い順にソートし、ソートされた構文解析木をRNA配列の二次構造の候補として出力するので、一配列に対して多文法で構文解析を行うことができるようになる。すなわち、各生成文法に対してそれぞれ構文解析し適合度計算を行い適合度を得る。その結果、生成文法ごとに適合度を得られることになり、それらの適合度をソートすることによって生成文法に順位を付ける。これにより、生成文法に対する構造トポロジーにも順位が付けられることになるので、最終的にRNA配列
- 10 列が取り得る可能性の高い順に構造トポロジーを確認することができるようになる。

- つぎの発明にかかるプログラムは、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、
- 20 上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力する出力ステップとを含むことを特徴とする。

- 25 このプログラムによれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、

計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力するので、多配列に対して一文法で構文解析を行うことができるようになる。すなわち、与えられた特定の構造トポロジーに対し、対応する生成文法を取得し、これをを用いてRNA配列データベースに格納されているすべてまたは一部のRNA配列をそれぞれ構文解析し、ある閾値以下の適合度で構文解析に成功したRNA配列群を結果として出力する。これにより、与えられた特定の二次構造を取り得るようなRNA配列を検索することができるようになる。

つぎの発明にかかるプログラムは、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出ステップと、上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出ステップにて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分をマークすることにより、上記RNA配列間で共通に有する構造トポロジーを可視化する共通構造マトリックス作成ステップとを含むことを特徴とする。

このプログラムによれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を抽出し、構造トポロジーとRNA配列とを二次元マトリックスで表示し、二次元マトリックスにおいて抽出されたRNA配列と構造トポロジーに対応する格子部分をマークすることにより、RNA配列間で共通に有する構造トポロジーを可視化するので、RNA配列間の共通構造を容易に発見することが

できるようになる。

つぎの発明にかかるプログラムは、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、利用者が入力したDNA配列から転写されるRNA配列を作成するRNA配列作成ステップと、上記RNA配列作成ステップにより作成された上記RNA配列に対して上記生成文法を適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列に対応する上記DNA配列部分を遺伝子の候補として予測する遺伝子予測ステップとを含むことを特徴とする。

このプログラムによれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、利用者が入力したDNA配列から転写されるRNA配列を作成し、作成されたRNA配列に対して生成文法を適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列に対応するDNA配列部分を遺伝子の候補として予測するので、DNA配列のうち既知のトポロジーを有する可能性のあるRNA配列に対応する部分について、遺伝子部分である可能性があることを予測することができるようになる。

つぎの発明にかかるプログラムは、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度に基づいて上記RNA配列間の類似度を計算する類似度計算ステップとを含むことを特徴とする。

このプログラムによれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度に基づいてRNA配列間の類似度を計算するので、RNA構造の類似度を容易に求めることができるようになる。

つぎの発明にかかるプログラムは、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出ステップと、上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出ステップにて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分に上記適合度を表示する適合度マトリックスを作成する適合度マトリックス作成ステップと、上記適合度マトリックス作成ステップにて作成された上記適合度マトリックスについて、上記適合度により上記構造トポロジーをソートし、他のRNA配列について当該ソートされた上記構造トポロジーの順番に対応する上記生成文法により構文解析を行い上記適合度が最大となる上記構文解析木を求め、予め定めた条件を満たす上記適合度を持つ上記構文解析木に対応する上記他のRNA配列を抽出する共通構造抽出ステップとを含むことを特徴とする。

このプログラムによれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を抽出し、構造トポロジーとRNA配列とを二次元マトリックスで表示し、二次元マトリックスにおいて抽出されたRNA配列と構造トポロジー

に対応する格子部分に適合度を表示する適合度マトリックスを作成し、適合度マトリックスについて、適合度により構造トポロジをソートし、他のRNA配列について当該ソートされた構造トポロジの順番に対応する生成文法により構文解析を行い適合度が最大となる構文解析木を求め、予め定めた条件を満たす適合度を持つ構文解析木に対応する他のRNA配列を抽出するので、共通の構造を持つRNA配列を容易に発見することができるようになる。

また、本発明は記録媒体に関するものであり、本発明にかかる記録媒体は、上記に記載されたプログラムを記録したことを特徴とする。

この記録媒体によれば、当該記録媒体に記録されたプログラムをコンピュータに読み取らせて実行することによって、上記に記載されたプログラムをコンピュータを利用して実現することができ、これら各プログラムと同様の効果を得ることができる。

図面の簡単な説明

第1図は、RNAの構造トポロジの一例を説明する図であり、第2図は、従来技術によるRNAの二次構造がステムループをとる場合の構文解析木の一例を示す図であり、第3図は、従来技術によるRNA二次構造の構造トポロジについて、文法が固定されるとそれに対応した構造トポロジが規定されることを示す概念図であり、第4図は、従来技術である本文法パーザを用いて、既知の文法からRNA配列の構文解析木を導出する場合の一例を示す図であり、第5図は、本発明が適用される本システムの構成の一例を示すブロック図であり、第6図は、文法データベース106bに格納される情報の一例を示す図であり、第7図は、本実施の形態における本システムのRNA二次構造予測処理の一例を示す処理概念図であり、第8図は、本実施の形態における本システムの同一構造RNA配列抽出処理の一例を示す処理概念図であり、第9図は、本実施の形態における本システムの共通構造抽出処理の一例を示す処理概念図であり、第10図は、本実施の形態における本システムの構造類似度計算処理の一例を示す処理概念図であり、

第 1 1 図は、本実施の形態における本システムの遺伝子予測処理の一例を示す処理概念図であり、第 1 2 図は、ペナルティ P と類似度ベクトル s_1 、 s_2 の概念を説明する図であり、第 1 3 図は、RNA 二次構造トポロジーの例を示す図であり、第 1 4 図は、 s_1 の構文解析木と二次構造を示す図であり、第 1 5 図は、塩基対の自由エネルギーを示す図であり、第 1 6 図は、ループの自由エネルギーを示す図であり、第 1 7 図は、それぞれの文法について $-\Delta G$ の適合度指標において最適な構文解析木とそれに対応する二次構造を示す図であり、第 1 8 図は、選択されたトポロジー集合のなかで s_2 が適合する構造候補を示す図であり、第 1 9 図は、選択されたトポロジーをとりうる配列の候補を示す図であり、第 2 0 図は、構文解析木の適合度を要素に持つマトリックスを示す図であり、第 2 1 図は、 s の最適な二次構造を示す図であり、第 2 2 図は、構文解析木の適合度を要素に持つマトリックスを示す図であり、第 2 3 図は、出力結果の一例を示す図である。

発明を実施するための最良の形態

以下に、本発明にかかる RNA 配列解析装置、RNA 配列解析方法、プログラム、および、記録媒体の実施の形態を図面に基づいて詳細に説明する。なお、この実施の形態によりこの発明が限定されるものではない。

特に、以下の実施の形態においては、本発明を、木文法に適用した例について説明するが、この場合に限られず、全ての生成文法において、同様に適用することができ。

[本システムの概要]

以下、本システムの概要について説明し、その後、本システムの構成および処理等について詳細に説明する。

このシステムは、概略的に、以下の基本的特徴を有する。すなわち、本システムの RNA 配列解析装置は、RNA 二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA 配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度を計算し、計算

された適合度のうち予め定めた条件を満たす適合度である構文解析木を適合度が高い順にソートし、ソートされた構文解析木をRNA配列の二次構造の候補として出力する。ここで、生成文法は、木文法、文脈自由文法などを含むが、シュエードノットを表現するためには木文法が最も適しているため、木文法を用いることが好ましい。

また、本装置は、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力する。

また、本装置は、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を抽出し、構造トポロジーとRNA配列とを二次元マトリックスで表示し、二次元マトリックスにおいて抽出されたRNA配列と構造トポロジーに対応する格子部分をマークすることにより、RNA配列間で共通に有する構造トポロジーを可視化する。

また、本装置は、利用者が入力したDNA配列から転写されるRNA配列を作成し、作成されたRNA配列に対して生成文法を適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列に対応するDNA配列部分を遺伝子の候補として予測する。

さらに、本装置は、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度に基づいてRNA配列間の類似度を計算する。

[システム構成]

まず、本システムの構成について説明する。第5図は、本発明が適用される本システムの構成の一例を示すブロック図であり、該構成のうち本発明に係る部分のみを概念的に示している。本システムは、概略的に、配列情報を解析する

RNA配列解析装置であるRNA配列解析装置100と、配列情報等に関する外部データベースやホモロジー検索用の外部分析プログラム等を提供する外部システム200とを、ネットワーク300を介して通信可能に接続して構成されている。

5 第5図においてネットワーク300は、RNA配列解析装置100と外部システム200とを相互に接続する機能を有し、例えば、インターネット等である。

第5図において外部システム200は、ネットワーク300を介して、RNA配列解析装置100と相互に接続され、利用者に対して配列情報等に関する外部データベースやホモロジー検索やモチーフ検索等の外部分析プログラムを実行する
10 ウェブサイトを提供する機能を有する。

ここで、外部システム200は、WEBサーバやASPサーバ等として構成してもよく、そのハードウェア構成は、一般に市販されるワークステーション、パーソナルコンピュータ等の情報処理装置およびその付属装置により構成してもよい。また、外部システム200の各機能は、外部システム200のハードウェア
15 構成中のCPU、ディスク装置、メモリ装置、入力装置、出力装置、通信制御装置等およびそれらを制御するプログラム等により実現される。

第5図においてRNA配列解析装置100は、概略的に、RNA配列解析装置100の全体を統括的に制御するCPU等の制御部102、通信回線等に接続されるルータ等の通信装置（図示せず）に接続される通信制御インターフェース部
20 104、入力装置112および出力装置114に接続される入出力制御インターフェース部108、および、各種のデータベースやテーブル（RNA配列データベース106a～共通構造マトリックス106c）を格納する記憶部106を備えて構成されており、これら各部は任意の通信路を介して通信可能に接続されている。さらに、このRNA配列解析装置100は、ルータ等の通信装置および専用線等の有線または無線の通信回線を介して、ネットワーク300に通信可能に
25 接続されている。

記憶部106に格納される各種のデータベース（RNA配列データベース10

6 a ~ 共通構造マトリックス 1 0 6 c) は、固定ディスク装置等のストレージ手段であり、各種処理に用いる各種のプログラムやテーブルやファイルやデータベースやウェブページ用ファイル等を格納する。

5 これら記憶部 1 0 6 の各構成要素のうち、RNA 配列データベース 1 0 6 a は、RNA 配列を格納したデータベースである。RNA 配列データベース 1 0 6 a は、インターネットを経由してアクセスする外部の RNA 配列データベースであってもよく、また、これらのデータベースをコピーしたり、オリジナルの配列情報を格納したり、さらに独自のアノテーション情報等を付加したりして作成したイン
10 ハウスデータベースであってもよい。また、RNA 配列データベース 1 0 6 a は、c DNA 等の DNA 配列データベースに基づいて予め生成された、あるいは必要時に動的に生成された RNA 配列を格納したのもでもよい。

また、文法データベース 1 0 6 b は、RNA 二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段である。ここで、第 6 図は、文法データベース 1 0 6 b に格納される情報の一例を示
15 す図である。第 6 図に示すように、文法データベース 1 0 6 b は、構造トポロジーと、その構造トポロジーに対応する文法とを対応付けて格納する。ここで、文法データベース 1 0 6 b には、第 6 図に示したように、構造トポロジーと文法とが 1 対 1 で対応するようにしてもよく、また、複数のトポロジーが結合した文法
(例えば、シュードノットとヘアピンループとを両方持つトポロジーなど) や、
20 特徴的な構造を有する RNA 用の文法 (例えば、r RNA に特有の構造トポロジーなど) や、所定のカテゴリーの RNA が共通で備えるトポロジー用の文法や、全ての RNA に適合する文法などを規定してもよい。

また、共通構造マトリックス 1 0 6 c は、構造トポロジーと RNA 配列とを二次元マトリックスで表示するためのテーブル (記憶領域) である。

25 また、第 5 図において、通信制御インターフェース部 1 0 4 は、RNA 配列解析装置 1 0 0 とネットワーク 3 0 0 (またはルータ等の通信装置) との間における通信制御を行う。すなわち、通信制御インターフェース部 1 0 4 は、他の端末

と通信回線を介してデータを通信する機能を有する。

また、第5図において、入出力制御インターフェース部108は、入力装置112や出力装置114の制御を行う。ここで、出力装置114としては、モニタ（家庭用テレビを含む）の他、スピーカを用いることができる（なお、以下において5 出力装置をモニタとして記載する）。また、入力装置112としては、キーボード、マウス、および、マイク等を用いることができる。また、モニタも、マウスと協働してポインティングデバイス機能を実現する。

また、第5図において、制御部102は、OS（Operating System）等の制御プログラム、各種の処理手順等を規定したプログラム、および10 所要データを格納するための内部メモリを有し、これらのプログラム等により、種々の処理を実行するための情報処理を行う。制御部102は、機能概念的に、構造予測部102a、類似度計算部102d、共通構造マトリックス作成部102f、および、遺伝子予測部102gを備えて構成されている。

このうち、構造予測部102aは、入力された既知の文法に従ってRNA配列15 の構文解析を行い構文解析木を導出する機能（構文解析部102b）、および、導出された構文解析木に対して適合度の計算を行う機能（適合度計算部102c）などを有する。

また、類似度計算部102dは、複数のRNA配列間の類似度を計算する類似度計算手段である。

20 また、共通構造マトリックス作成部102fは、適合度計算手段により計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を抽出する抽出手段、および、構造トポロジーとRNA配列とを二次元マトリックスで表示し、二次元マトリックスにおいて抽出手段にて抽出されたRNA配列と構造トポロジーに対応する格子部分をマークすることにより、RNA25 配列間で共通に有する構造トポロジーを可視化する共通構造マトリックス作成手段、二次元マトリックスにおいて抽出手段にて抽出されたRNA配列と構造トポロジーに対応する格子部分に適合度を表示する適合度マトリックスを作成する適

合度マトリックス作成手段、および、適合度マトリックス作成手段にて作成された適合度マトリックスについて、適合度により構造トポロジーをソートし、他のRNA配列について当該ソートされた構造トポロジーの順番に対応する生成文法により構文解析を行い適合度が最大となる構文解析木を求め、予め定めた条件を
5 満たす適合度を持つ構文解析木に対応する他のRNA配列を抽出する共通構造抽出手段である。

また、遺伝子予測部102gは、利用者が入力したDNA配列から転写されるRNA配列を作成するRNA配列作成手段、および、適合度計算手段により計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出した
10 RNA配列に対応するDNA配列部分を遺伝子の候補として予測する遺伝子予測手段である。なお、これら各部によって行なわれる処理の詳細については、後述する。

[システムの処理]

次に、このように構成された本実施の形態における本システムの処理の一例について、以下に第7図～第11図を参照して詳細に説明する。
15

[RNA二次構造予測処理]

まず、RNA二次構造予測処理の詳細について第7図を参照して説明する。第7図は、本実施の形態における本システムのRNA二次構造予測処理の一例を示す処理概念図である。

20 まず、文法データベース106bに既知のRNAの構造トポロジーを表す文法を集積する。そして、利用者が構造未知のRNA配列であってその二次構造を特定したいものを入力装置112を介してRNA配列解析装置100に入力すると（ステップSA-1）、構造予測部102aは、構文解析部102bの処理により、文法データベース106bから文法を取り出し（ステップSA-2）、RNA
25 A配列に対して各文法を適合して構文解析（パース）を行う（ステップSA-3）。ここで、利用者のRNA配列の入力は、RNA配列データベース106aから所望の配列を選択することにより入力してもよく、外部システム200の外部

データベースから所望の配列を選択することにより入力してもよく、さらに、所望の配列を直接入力してもよい。

そして、構造予測部 102 a は、適合度計算部 102 c の処理により、パーズが成功し導出された構文解析木について、例えば、ループや、塩基対とその他の二次構造要素のそれぞれの自由エネルギーの合計などを計算することにより求める平衡自由エネルギー (ΔG) などに基づいて適合度を求める。ここで、適合度計算方法は、例えば上述した文献 1 から 3 に示した方法のほか、従来のいずれの方法を用いてもよい。

そして、構造予測部 102 a は、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を適合度が高い順にソートする（ステップ S A-4）。

そして、構造予測部 102 a は、入出力制御インターフェース部 108 を介して出力装置 114 にソートされた構文解析木とその適合度を出力することにより、利用者が入力した一配列に対して多文法で構文解析を行うことができるようになる。すなわち、各生成文法に対してそれぞれ構文解析し適合度計算を行い適合度を得る。その結果、生成文法ごとに適合度を得られることになり、それらの適合度をソートすることによって生成文法に順位を付ける。これにより、生成文法に対する構造トポロジーにも順位が付けられることになるので、最終的に RNA 配列が取り得る可能性の高い順に構造トポロジーを確認することができるようになる。これにて、RNA 二次構造予測処理が終了する。

[同一構造 RNA 配列抽出処理]

次に、同一構造 RNA 配列抽出処理の詳細について第 8 図を参照して説明する。第 8 図は、本実施の形態における本システムの同一構造 RNA 配列抽出処理の一例を示す処理概念図である。

まず、利用者は、特定の構造トポロジーに対応する文法を文法データベース 106 b から選択する。そして、構造予測部 102 a は、構文解析部 102 b の処理により、RNA 配列データベース 106 a から RNA 配列を取り出し（ステッ

プSB-1)、各RNA配列に対して文法を適合して(ステップSB-2)、構文解析を行う(ステップSB-3)。

そして、適合度計算部102cは、導出された構文解析木に対して適合度計算を行う。そして、構造予測部102aは、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を、指定した文法が表す構造トポロジーの二次構造を持つRNA配列の候補として抽出する(ステップSB-4)。

そして、構造予測部102aは、抽出されたRNA配列を当該文法が表す構造トポロジーの二次構造をもつ可能性のあるものとして、入出力制御インターフェース部108を介して出力装置114に出力する(ステップSB-5)。これにて、同一構造RNA配列抽出処理が終了する。

[共通構造抽出処理]

次に、共通構造抽出処理の詳細について第9図を参照して説明する。第9図は、本実施の形態における本システムの共通構造抽出処理の一例を示す処理概念図である。

まず、構造予測部102aは、RNA配列データベース106aから1つまたは2つ以上のRNA配列を取り出し(ステップSC-1およびステップSC-2)、構文解析部102bの処理により、各RNA配列に対して、文法データベース106bから取り出した(ステップSC-3)、1つまたは2つ以上の文法を適合する(ステップSC-4)。RNA配列解析装置100は、これらの取り出しや、パース処理について並列処理を行ってもよく、また、順次処理を行ってもよい。

そして、適合度計算部102cは、導出された構文解析木に対して適合度計算を行い、共通構造マトリックス作成部102fの処理により、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を抽出する(ステップSC-5)。

そして、共通構造マトリックス作成部102fは、適合した文法が表す構造ト

ポロジとRNA配列とを二次元マトリックスで表示し、二次元マトリックスにおいて抽出されたRNA配列と構造トポロジに対応する格子部分をマークすることにより、RNA配列間で共通に有する構造トポロジを可視化する（ステップSC-6）。

- 5 ここで、マークは、第9図に示すように対象の格子部分に特定の色をつけてもよく、また、特定の記号（○など）や文字（「Y」など）を記載してもよい。これにより、例えば、縦方向にマークが連続した場合（第9図に示す例においては、2番目の構造トポロジの列）には、この構造トポロジが各RNA配列に共通に持っている配列であることが可視的に把握することができるようになる。これ
- 10 にて、共通構造抽出処理が終了する。

[構造類似度計算処理]

次に、構造類似度計算処理の詳細について第10図を参照して説明する。第10図は、本実施の形態における本システムの構造類似度計算処理の一例を示す処理概念図である。

- 15 まず、利用者が、入力装置112により類似度を計算したい複数（第10図の例では2個）のRNA配列をRNA配列解析装置100に入力する（ステップSE-1）。

- そして、類似度計算部102dは、文法データベース106bから1つまたは2つ以上の文法を取り出し（ステップSE-2）、構文解析部102bの処理により、入力したRNA配列について文法を適合して構文解析を行う（ステップSE-3）。また、適合度計算部102cは、導出された構文解析木に対して適合度計算を行う（ステップSE-4）。
- 20

- そして、類似度計算部102dは、文法を適合することにより導出された構文解析木と適合度（導出されなかった場合には、それを表現する特別な値を設定する）を各RNA配列ごとに対応付けてベクトル演算や内積の計算などを行うことにより（ステップSE-5）、RNA配列間の類似度を計算する（ステップSE-6）。
- 25

例えば、入力である i 個の RNA 配列を $RNA_1, RNA_2, \dots, RNA_i$ とし、文法データベース 106b に格納されている N 個の文法を G_1, G_2, \dots, G_N とし、また、RNA 配列 x と文法 g に対するパーザが成功したときの適合度を $r(x, g)$ とする。ここで、適合度は実数値とし、適合度が大きいほどその構造を取り易いことを示すものとする。

また、入力 RNA_j に関する適合度のベクトル R_j において、 R_j の k 番目の要素 $R_j[k]$ は、 RNA_j と G_k に対するパーザが成功したときは $r(RNA_j, G_k)$ とし、 RNA_j と G_k に対するパーザが失敗したときは仮に「×」とする。

このとき、類似度計算部 102d による類似度計算は、以下の手法により行われる。まず、2つの RNA 配列の適合度のベクトル R_1 と R_2 を入力する。

ついで、類似度計算部 102d は、類似度ベクトル S_1, S_2 とペナルティ P を求める。ここで、「ペナルティ P 」は、 $R_1[k]$ と $R_2[k]$ の片方だけが「パーザ失敗(×)」である k の個数であり、「類似度ベクトル S_1, S_2 」は、 $R_1[k]$ も $R_2[k]$ も「パーザ失敗(×)」ではない箇所だけを抜き出したベクトルである。第 12 図は、ペナルティ P と類似度ベクトル S_1, S_2 の概念を説明する図である。

ついで、類似度計算部 102d は、類似度ベクトル S_1, S_2 間の距離 D を以下の方法により求める。まず、類似度ベクトル S_1, S_2 の要素数(ベクトルの次元)を M とする。そして、類似度計算で一般的に用いられるユークリッド距離を用いて以下の数式により距離を計算する。

$$D = \sqrt{\sum \{ (S_1[k] - S_2[k])^2 \}}$$

($\sqrt{\quad}$ は平方根であり、 \sum は $k = 1 \sim M$ に関する総和である。)

ここで、距離 D が大きい場合には類似度が低いことになり、また、ペナルティ P が大きいと類似度が低いことになるので、ペナルティ P と距離 D を用いて以下の数式により類似度 Sim を求める。

$$Sim = a^P / D$$

(a は定数 ($0 < a < 1$) である。)

- 5 そして、 Sim を類似度として出力する。ここで、定数 a を小さくすると、距離 D よりもペナルティ P が重視されることになる。これにて、構造類似度計算処理が終了する。

[遺伝子予測処理]

- 10 次に、遺伝子予測処理の詳細について第11図を参照して説明する。第11図は、本実施の形態における本システムの遺伝子予測処理の一例を示す処理概念図である。

- 15 まず、利用者が遺伝子部分が未知のDNA配列を入力装置112を介してRNA配列解析装置100に入力すると、RNA配列解析装置100は、遺伝子予測部102gの処理により、入力されたDNA配列に基づいて、当該DNA配列から転写されるRNA配列（以下、「予測RNA配列」という）を自動的に変換して作成する（ステップSF-1）。ここで、利用者のDNA配列の入力は、外部システム200の外部データベースやインハウスデータベースから所望のDNA配列を選択することにより入力してもよく、さらに、所望の配列を直接入力してもよい。

- 20 ついで、構造予測部102aがこの予測RNA配列を構文解析部102bに入力すると（ステップSF-2）、構文解析部102bの処理により、文法データベース106bから1つまたは2つ以上の文法が取り出され（ステップSF-3）、各文法を予測RNA配列に適合する（ステップSF-4）。

- 25 そして、適合度計算部102cは、構文解析部102bが導出した構文解析木について適合度計算を行い（ステップSF-5）、遺伝子予測部102gは、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出した予測RNA配列に対応するDNA配列部分を遺伝子の候補として予測する（ス

テップSF-6)。すなわち、DNA配列のうち、当該予測RNA配列の部分が遺伝子部分である可能性が高い領域として出力される。

- これにより、DNA配列のうち既知のトポロジーを有する可能性のある予測RNA配列に対応する部分について、遺伝子部分である可能性があることを予測することができるようになる。これにて、遺伝子予測処理が終了する。

[実施例]

本発明の実施例について、以下に第13図～第23図を参照して説明する。

1 準備

- 本節では実施例の準備として、いくつかの具体的なRNA二次構造トポロジーを定義し、それらをモデリングする生成文法を規定する。本実施例では説明の便宜上、生成文法として文脈自由文法を用いるが、よりモデリング能力の高いRNA木文法（文献1）を用いた場合でも同様のことが説明できる。

1. 1 二次構造トポロジー

第13図に示される2つのRNA二次構造トポロジーを考える。

- ステムループは、ステム($H(a)$)とヘアピンループ($L(a)$)から構成される。2並列ステムループは、並列に並んだ2つのステムループから構成される。それぞれのステム部分($H_1(b)$ 、 $H_2(b)$)とヘアピンループ部分、($L_1(b)$ 、 $L_2(b)$)の他にステムとステムをつなぐループ部分($I(b)$)がある。

- 上記構造トポロジーについて、さらに具体的な特徴を考えることができる。例えばステムやループ部分のサイズの制約、ステムを構成する塩基対にミスマッチ（内部ループやバルジループ）を許すかどうか、あるいは特定の場所に特定の塩基配列を含むかどうか、などといったより詳細な特徴を持ったトポロジーを考えることが可能である。そこで、本実施例では次のような特徴を持ったRNA二次構造トポロジー T_1 、 T_2 を扱う。

トポロジー T_1

- ー 以下の特徴を持ったステムループ構造（第13図(a)参照）である。

- ステム (H (a)) を構成する塩基対はミスマッチを含まない。
- ステム (H (a)) サイズは1塩基対以上とする。
- ヘアピンループ (L (a)) のサイズは1塩基以上とする。

トポロジー T_2

- 5 — 以下の特徴を持った2並列ステムループ構造 (第13図 (b) 参照) である。
- トポロジー T_1 を2つ並列に並べたもの。
 - ステム ($H_1 (b)$) とステム ($H_2 (b)$) の間のループ (I (b)) の長さは1塩基以上とする。

10 1. 2 文脈自由文法による二次構造トポロジーのモデリング

以上のように定義された2つのトポロジー T_1 、 T_2 を文脈自由文法を用いてモデリングする。文脈自由文法は一般に以下の4項組によって規定される。

$$G = (N, \Sigma, P, S)$$

15

Nは非終端記号の有限集合、 Σ は終端記号の有限集合、Pは生成規則の有限集合、Sは開始記号を表す。

しかしながら、本実施例では常に $\Sigma = \{a, u, g, c\}$ 、開始記号はS、さらにNは生成規則Pに出現する非終端記号のみからなるものとすることから、

- 20 Pのみを指定することにより文脈自由文法Gを規定することができる。よって便宜上、本稿では文脈自由文法Gを規定する際には、生成規則の有限集合Pのみを指定することにする。

(1) トポロジー T_1 は以下の生成規則からなる文脈自由文法 G_1 によってモデリングされる。

25

$$\begin{aligned}
S &\rightarrow xH\bar{x} \\
H &\rightarrow xH\bar{x} \mid L \\
L &\rightarrow xL \mid x
\end{aligned}$$

5

ただし、 $x \in \Sigma$ とし、 \bar{x} は x と塩基対を形成する相手の塩基とする。

すなわち、ワトソン・クリック塩基対のみを考える場合は、一番最初の生成規則は以下と同値である。

$$10 \quad S \rightarrow aHu \mid uHa \mid gHc \mid cHg$$

非ワトソン・クリック塩基対を許す場合はさらに、 $S \rightarrow gHu$ などを追加してもよい。

G_1 において、

15

$$S \rightarrow xH\bar{x} \text{ と } H \rightarrow xH\bar{x}$$

によって塩基対（ステムを構成）が生成され、 $L \rightarrow xL$ と $L \rightarrow x$ によって塩基対を形成しない塩基（ループを構成）が生成されるとみなす。すると、 G_1 は RNA の二次構造を生成することができることになる。このように、任意の文脈自由文法 G に対して、 G が生成することができるすべての RNA 二次構造からなる集合 $SS(G)$ が規定されることになる。

G_1 がトポロジー T_1 をモデリングする、とは以下が成り立つことを言う。“ G_1 はトポロジー T_1 の性質を満たすすべての RNA 二次構造を生成することができ、
25 なおかつ、 G_1 で生成することができるすべての RNA 二次構造はトポロジー T_1 の性質を満たす。”

これは、 G_1 による導出をみれば自明である。 G_1 による導出はすべて以下のよ

うになる。

ただし、 $n \geq 1$ 、 $l \geq 1$ とする。

$$\begin{aligned}
 5 \quad S &\rightarrow x_1 H \overline{x_1} \\
 &\rightarrow x_1 x_2 H \overline{x_2 x_1} \rightarrow \dots \rightarrow x_1 x_2 \dots x_n H \overline{x_n \dots x_2 x_1} \\
 &\rightarrow x_1 x_2 \dots x_n L \overline{x_n \dots x_2 x_1} \\
 &\rightarrow x_1 x_2 \dots x_n y_1 L \overline{x_n \dots x_2 x_1} \\
 &\rightarrow \dots \rightarrow x_1 x_2 \dots x_n y_1 \dots y_{l-1} L \overline{x_n \dots x_2 x_1} \\
 &\rightarrow x_1 x_2 \dots x_n y_1 \dots y_{l-1} y_l \overline{x_n \dots x_2 x_1}
 \end{aligned}$$

10

このとき、 $x_1 x_2 \dots x_n$ と、これに対応する $\overline{x_n \dots x_2 x_1}$ の部分がステムに、 $y_1 \dots y_{l-1} y_l$ の部分がヘアピンループに対応する。また、 $n \geq 1$ 、 $l \geq 1$ であるからステムのサイズは1塩基対以上、ヘアピンループのサイズは1塩基以上である。

15 よって、 G_1 は T_1 をモデリングすることがわかる。

(2) トポロジー T_2 は以下の生成規則からなる文脈自由文法 G_2 によってモデリングされる。

$$\begin{aligned}
 20 \quad S &\rightarrow S_1 I S_2 \\
 S_1 &\rightarrow x H \bar{x} \\
 S_2 &\rightarrow x H \bar{x} \\
 H &\rightarrow x H \bar{x} \mid L \\
 L &\rightarrow x L \mid x \\
 I &\rightarrow x L
 \end{aligned}$$

25 以下の生成規則からなる文脈自由文法 G_0 は、文脈自由文法によって生成することが可能なすべてのRNA二次構造を生成することができる万能な文脈自由文法である。

$$S \rightarrow SS | xS\bar{x} | xS | Sx | x | \lambda$$

ただし、 λ は空文字を表す。例えば、 G_1 によるいかなる導出も G_0 によって
 5 シミュレートできる。すなわち、以下のような導出を G_0 によって行なうことが
 可能である。

$$\begin{aligned} S &\rightarrow x_1 S \bar{x}_1 \\ &\rightarrow x_1 x_2 S \bar{x}_2 \bar{x}_1 \rightarrow \dots \rightarrow x_1 x_2 \dots x_n S \bar{x}_n \dots \bar{x}_2 \bar{x}_1 \\ &\rightarrow x_1 x_2 \dots x_n S \bar{x}_n \dots \bar{x}_2 \bar{x}_1 \\ 10 &\rightarrow x_1 x_2 \dots x_n y_1 S \bar{x}_n \dots \bar{x}_2 \bar{x}_1 \\ &\rightarrow \dots \rightarrow x_1 x_2 \dots x_n y_1 \dots y_{l-1} S \bar{x}_n \dots \bar{x}_2 \bar{x}_1 \\ &\rightarrow x_1 x_2 \dots x_n y_1 \dots y_{l-1} y_l \bar{x}_n \dots \bar{x}_2 \bar{x}_1 \end{aligned}$$

15 上記の導出は非終端記号以外、すなわち生成される RNA 二次構造は、 G_1 に
 よって生成されるものとまったく同じである。よって G_1 が生成可能なすべての
 二次構造を G_0 によって生成可能なことがわかる。すなわち、

$$SS(G_0) \supseteq SS(G_1)$$

20

である。

このように、どんな文脈自由文法 G に対しても

$$SS(G_0) \supseteq SS(G)$$

25

が成り立つことが知られている。以降では、このような万能文法によって生成さ
 れる二次構造全体を“すべての二次構造”と考える。

1. 3 構文解析木と適合度

ある与えられたRNA配列がある与えられたRNA二次構造トポロジーの性質を満たす二次構造を形成できるかどうかという問題は、対象トポロジーをモデリングした文法によって対象配列が導出できるかという問題に対応する。これは生成文法の構文解析アルゴリズムによって解くことができる。

構文解析アルゴリズムは、与えられた文法によって、与えられた配列が導出できるかどうかを判定し、導出可能な場合はその導出過程、すなわち構文解析木を出力する。二次構造トポロジーをモデリングした文法において、構文解析木は二次構造を表現しているので、構文解析アルゴリズムは、対象トポロジーに適合する具体的な二次構造を出力すると解釈してよいことになるからである。

RNA配列 $s_1 = g g g g a a a c c c c$ がトポロジー T_1 、 T_2 に適合する二次構造を形成できるかどうかについて考える。

配列 s_1 は G_1 によって以下のように導出できる。これにより配列 s_1 は T_1 に適合する二次構造をとりうることがわかる。

$$\begin{aligned} S &\rightarrow g H c \rightarrow g g H c c \rightarrow g g g H c c c \rightarrow g g g g H c c c c \\ &\rightarrow g g g g L c c c c \rightarrow g g g g a L c c c c \rightarrow g g g g a a L c c c c \\ &\rightarrow g g g g a a a c c c c \end{aligned} \quad (1)$$

また、 s_1 は G_1 によって以下のようにも導出できる。

$$\begin{aligned} S &\rightarrow g H c \rightarrow g g H c c \rightarrow g g g H c c c \\ &\rightarrow g g g L c c c \rightarrow g g g g L c c c \rightarrow g g g g a L c c c \\ &\rightarrow g g g g a a L c c c \rightarrow g g g g a a a L c c c \\ &\rightarrow g g g g a a a c c c c \end{aligned} \quad (2)$$

しかし、 s_1 は G_2 によって導出することはできない。これにより s_1 はトポロ

ジー T_2 に適合する二次構造をとりえないことがわかる。

s_1 を G_1 によって上記の2通りの方法で導出するとき、それぞれの導出に対応する構文解析木とそれが表現する二次構造を第14図に示す。すなわち、(1)のように導出した場合は、第14図の(1)に示される構文解析木と二次構造が生成され、(2)のように導出した場合は、第14図の(2)に示される構文解析木と二次構造が生成される。

この例のように複数の構文解析木が得られたときに、どの構文解析木、すなわち、どの二次構造を結果として出力するかを決定する必要がある。そのために、何らかの評価関数によって構文解析木（あるいは二次構造）にスコアを与え、構文解析木（あるいは二次構造）に順位を付ける必要がある。このようなスコアは文法によって異なる評価関数を用いても良いし、文法に依存しない絶対的な評価関数であってもよい。以降では、このスコアを適合度と呼ぶ。

以下に、これまでに利用されている適合度の評価法の例を示すが、本発明によって利用される適合度は以下のものに限定されない。

15 (1) 塩基対数による適合度の評価

一般に、塩基対を形成する際の水素結合によってRNA分子はエネルギー的に安定なものになる。そこでこの評価法では、単純に塩基対の数が多い二次構造ほど優先させる。つまり、構文解析木の適合度として、対応する二次構造の塩基対数を用いる。この評価法において、上記の例の適合度を評価すると、第14図の(1)に示される構文解析木は適合度3となり、(2)に示される構文解析木は適合度2となり、適合度の大きい(1)の構造が採用されることになる。

本評価法に基づいた代表的な手法として、Nussinovの折り畳みアルゴリズム [Nussinov, R., Pieczenik, G., Geiggs, J. R., and Kleitman, D. J., "Algorithms for loop matchings," SIAM journal of Applied Mathematics, 35, 68-82, 1978] がある。

(2) 平衡自由エネルギー (ΔG) による適合度の評価

RNA 二次構造の物理化学的な安定度を計算するために、小さなモデル RNA 分子の熱力学的な実験によって決定された平衡自由エネルギー (ΔG) パラメータを利用する方法がある。ある二次構造の (ΔG) は、それを構成する塩基対やループなどの二次構造要素に対する自由エネルギーの合計で近似される。この自由エネルギーパラメータでは、塩基対により構造が安定化し、ループにより構造が不安定化する。各二次構造要素の詳細なパラメータは [Turner, D. H., Sugimoto, N., Jaeger, J. A., Longfellow, C. E., Freier, S. M., and Kierzek, R., "Improved parameters for prediction of RNA structure," Cold Spring Harbor Symposia Quantitative Biology, 52, 123-133, 1987] に示されている。ここでは塩基対の自由エネルギーを第 15 図に、ループの自由エネルギーを第 16 図に示す。

上記の自由エネルギーパラメータを用いて第 14 図の構造 (1) と (2) の (ΔG) を求めると、それぞれ以下ようになる。

$$\begin{aligned}\Delta G (\text{構造 (1)}) &= \Delta G (\text{gc, gc}) + \Delta G (\text{gc, gc}) \\ &\quad + \Delta G (\text{gc, gc}) \\ &\quad + (\Delta G) (\text{サイズ 3 のヘアピンループ}) \\ &= (-2.9) + (-2.9) + (-2.9) \\ &\quad + 7.4 = -1.3\end{aligned}$$

$$\begin{aligned}\Delta G (\text{構造 (2)}) &= \Delta G (\text{gc, gc}) + \Delta G (\text{gc, gc}) \\ &\quad + \Delta G (\text{サイズ 5 のヘアピンループ}) \\ &= (-2.9) + (-2.9) + 4.4 = -1.4\end{aligned}$$

ここで注意すべきことは、塩基対の自由エネルギーの計算方法である。連続して積み重なった 2 組の塩基対に対してひとつのエネルギー値が与えられる。すな

わち、構造（１）では５'側から数えて１番目の g c 塩基対と２番目の g c 塩基対に対して、 $\Delta G(g c, g c)$ が計算され、２番目の g c 塩基対と３番目の g c 塩基対に対して、 $\Delta G(g c, g c)$ が計算され、３番目の g c 塩基対と４番目の g c 塩基対に対して、 $\Delta G(g c, g c)$ が計算される。これに対し構造（
 5 2）では５'側から数えて１番目の g c 塩基対と２番目の g c 塩基対に対して、 $\Delta G(g c, g c)$ が計算され、２番目の g c 塩基対と３番目の g c 塩基対に対して、 $\Delta G(g c, g c)$ が計算される。

構文解析木の適合度を $-\Delta G$ と定めると、（１）の適合度は 1.3 となり、（
 10 2）の適合度は 1.4 となり、結果として適合度の大きい（２）の構造が採用されることになる。

ΔG に基づいた代表的な RNA 二次構造予測システムとして、Z u k e r の M f o l d（文献 3）がある。

（３） 導出確率による適合度の評価

確率文法とは個々の生成規則にその適用確率が付加された生成文法である。例
 15 えば G_1 の各生成規則に以下のような確率 p が付加されている確率文脈自由文法 G_1 を考える。

- $p(S \rightarrow a H u) = 0.2$
- $p(S \rightarrow u H a) = 0.2$
- 20 $p(S \rightarrow g H c) = 0.3$
- $p(S \rightarrow c H g) = 0.3$
- $p(H \rightarrow a H u) = 0.2$
- $p(H \rightarrow u H a) = 0.2$
- $p(H \rightarrow g H c) = 0.3$
- 25 $p(H \rightarrow c H g) = 0.2$
- $p(H \rightarrow L) = 0.1$
- $p(L \rightarrow a L) = 0.2$

$$p(L \rightarrow u L) = 0.2$$

$$p(L \rightarrow g L) = 0.15$$

$$p(L \rightarrow c L) = 0.15$$

$$p(L \rightarrow a) = 0.1$$

$$5 \quad p(L \rightarrow u) = 0.1$$

$$p(L \rightarrow g) = 0.05$$

$$p(L \rightarrow c) = 0.05$$

このとき G_1 による s_1 の導出確率は次のようにして計算される。すなわち、（

10 1) の導出確率は、

$$\begin{aligned} & p(S \rightarrow g H c) \times p(H \rightarrow g H c) \times p(H \rightarrow g H c) \times p(H \rightarrow g H c) \times \\ & p(H \rightarrow L) \times p(L \rightarrow a L) \times p(L \rightarrow a L) \times p(L \rightarrow a) \\ & = 0.3 \times 0.3 \times 0.3 \times 0.3 \times 0.1 \times 0.2 \times 0.2 \times 0.1 \end{aligned}$$

$$15 \quad = 0.00000324$$

と計算される。また、(2) の導出確率は、

$$\begin{aligned} & p(S \rightarrow g H c) \times p(H \rightarrow g H c) \times p(H \rightarrow g H c) \times p(H \rightarrow L) \times p(\\ 20 \quad & L \rightarrow g L) \times p(L \rightarrow a L) \times p(L \rightarrow a L) \times p(L \rightarrow a L) \times p(L \rightarrow c) \\ & = 0.3 \times 0.3 \times 0.3 \times 0.1 \times 0.15 \times 0.2 \times 0.2 \times 0.2 \\ & \quad \times 0.05 \end{aligned}$$

$$= 0.000000162$$

25 となる。

そこで構文解析木の適合度として導出確率の自然対数をとると、(1) の適合度は $1 \ln 0.00000324 = -12.6$ 、(2) の適合度は $1 \ln 0.000$

000162 = -15.6 となり、結果として適合度の大きい (1) の構造が採用されることになる。

本評価法の根本である、各生成規則に付加されるべき確率パラメータは、最尤推定法と内側・外側アルゴリズム (inside-outside algorithm) などにより学習してもよいし、ヒューリスティクス (heuristics) などによって主観的に推定してもよい。例えば、文献 [Sakakibara 5 "Stochastic Context-free Grammars for tRNA modeling," Nucleic Acids Research, 22, 5112-5120, 1994.] では複数の tRNA 10 A 配列から tRNA の構造をモデリングする確率文脈自由文法を学習する手法について述べられている。

以上いくつかの適合度評価法について説明したが、以降の説明では適合度として $-\Delta G$ を用いる。

次に RNA 配列 $s_2 = g c c c a u a g g c a a a g c c u a u g g g c$ がト 15 ポロジー T_1 、 T_2 に適合する二次構造を形成できるかどうかを考える。この場合も同様に s_2 が G_1 、 G_2 によって導出できるかどうかを調べればよい。結論からいうと、 s_2 は G_1 、 G_2 のいずれからでも導出することができる。さらにどちらの文法でも複数の導出のしかたが存在する。それぞれの文法について $-\Delta G$ の適合度指標において最適な構文解析木とそれに対応する二次構造を第 17 図に示す。

20 それぞれの構造の ΔG を求めると、以下のようになる。

$$\begin{aligned}
 \Delta G (\text{構造 (1)}) &= \Delta G (g c, c g) \times 2 + \Delta G (c g, c g) \\
 &\quad \times 2 + \Delta G (c g, a u) + \Delta G (a u, u a) \\
 &\quad + \Delta G (u a, a u) + \Delta G (a u, g c) \\
 &\quad + \Delta G (g c, g c) \\
 &\quad + \Delta G (\text{サイズ 3 のヘアピンループ}) \\
 &= (-3.4) \times 2 + (-2.9) \times 2 + (-1.8)
 \end{aligned}$$

$$\begin{aligned}
 &+ (-0.9) + (-1.1) + (-1.7) \\
 &+ (-2.9) + 7.4 = -13.6
 \end{aligned}$$

$$\begin{aligned}
 \Delta G (\text{構造 (2)}) &= \Delta G (g c, c g) \times 2 + \Delta G (c g, c g) \times 2 \\
 5 \quad &+ \Delta G (\text{サイズ4のヘアピンループ}) \times 2 \\
 &= (-3.4) \times 2 + (-2.9) \times 2 + 5.9 \times 2 \\
 &= -6.7
 \end{aligned}$$

よってトポロジー T_1 に適合するRNA二次構造において s_2 がとりうる最適な
 10 構造の適合度は13.6であることがわかる。また、トポロジー T_2 に適合する
 RNA二次構造において s_2 がとりうる最適な構造の適合度は6.7であること
 がわかる。また、 s_2 を万能な文法 G_0 によって構文解析すると、最適構造として
 構造(1)が見つかる。これにより、構造(1)が“すべての二次構造”の中で
 最適な構造であることがわかる。このように万能文法によって構文解析を行なう
 15 ことにより、すべての構造の中から最適な構造を見つけ出すことができる。

本発明の基本となる“RNA配列を生成文法に適用して構文解析木を導出する
 構文解析手段と、上記構文解析手段にて導出された構文解析木に対して適合度の
 計算を行う適合度計算手段と、適合度最大の構文解析木に対応する二次構造を出
 力する最適二次構造出力手段”は、適合度計算を組み込んだ構文解析アルゴリズム
 20 ムにより実現されることになる。このような構文解析アルゴリズムを構造予測ア
 ルゴリズムと呼ぶ。 ΔG の適合度を指標にしたRNA木文法に基づく構造予測ア
 ルゴリズムは文献1に示されている。

2. 発明実施例

本節では、前節までに定義したRNA配列 s_1 、 s_2 、トポロジー T_1 、 T_2 および
 25 それらをモデリングする文脈自由文法 G_0 、 G_1 、 G_2 、さらに適合度として $-\Delta G$
 を用いた場合の実施例について示す。

はじめに、“RNA二次構造の構造トポロジーと、当該構造トポロジーに適合

する生成文法とを対応付けて格納する文法格納手段”においては、例えば（L e
u - t R N A , G ' ） や （ 1 6 S r R N A , G " ） などのようにある構造トポ
ロジーに付けられた名称とその構造トポロジーをモデリングした文法が対応づけ
られて格納されている。本実施例では（ステムループ T_1 , G_1 ）と（2 並列ステ
ムループ T_2 , G_2 ）を含むような文法DBを仮定する。また、RNA配列 s_1 と s
5 s_2 を含むRNA配列DBを仮定する。

（1）生成文法と適合度計算で構造候補を出力

あるRNA配列に対して、これがとりうる構造トポロジーを適合度が大きい順
に知りたいとき、本発明によれば、以下の手順でこれを調べることができる。例
10 として、入力配列を s_2 、検査対象トポロジー集合を T_1 、 T_2 とした場合について
示す。

手順1） RNA配列を配列DBから指定する。あるいは直接入力する。ここで
は s_2 を指定する。

手順2） 検査対象とするトポロジーの集合（生成文法の集合）を文法DBから
15 選択する。ここでは T_1 と T_2 （ G_1 と G_2 ）を選択する。

手順3） 適合度の閾値を設定する。閾値は手順2で得られた各トポロジー（生
成文法）に対してそれぞれ設定してもよいし、共通にひとつ設定してもよい。こ
こでは T_1 （ G_1 ）に対して10を T_2 （ G_2 ）に対して5を設定する。

手順4） 手順1で得られた配列を手順2で得られた各生成文法によってそれぞ
20 れ構文解析を行ない、適合度最大の構文解析木を求める。ここでは s_2 を G_1 によ
って構文解析し、最大の適合度13.6を持つ構文解析木を得る（第17図（1）
参照）。

さらに s_2 を G_2 によって構文解析し、最大の適合度6.7を持つ構文解析木を
得る（第17図（2）参照）。

25 手順5） 手順4で得られた構文解析木のうち手順3で得られた閾値以上の適合
度をもつ構文解析木を適合度の大きい順にソートする。手順4で得られた G_1 に
対する適合度13.6の構文解析木1は手順3で G_1 に対して設定された閾値1.

0 よりも大きいのでこれをソート対象とする。手順4で得られた G_2 に対する適合度6.7の構文解析木2は手順3で G_2 に対して設定された閾値5より大きいのでこれをソート対象とする。上でソート対象となった構文解析木を適合度の大きい順にソートすることによって、構文解析木1、構文解析木2の順に順序づけされる。

手順6) 手順5でソートされた構文解析木の順に、対応するトポロジー名、適合度、構文解析木(二次構造)などを出力する。構文解析木1に対応して、ステムループ T_1 、適合度13.6、第17図(1)に示された二次構造を出力する。構文解析木2に対応して、2並列ステムループ T_2 、適合度6.7、第17図(2)に示された二次構造を出力する。

以上の結果から、選択されたトポロジー集合のなかで s_2 が適合する構造候補が第18図のように出力される。

従来の二次構造予測プログラムでは、与えられた配列がとりうる構造のなかで最適あるいは準最適な二次構造を順に出力するだけで、出力された構造がどのようなトポロジーであるかはユーザが判断しなければならなかった。本発明によれば、構造とトポロジーとを対応付けて出力することができるので、予測結果の確認にかかる労力が大きく軽減されることが期待される。

また、本発明の実施について厳密に上記の手順と同じである必要はない。例えば、手順1と2は順序をいれかえてもよいし、手順5の閾値により構文解析木を取捨選択する部分は手順4の構文解析部分に含めてしまってもよい。

(2) 同じ構造を持つ配列の候補を出力

ある構造トポロジーに対して、これに適合する二次構造をとりうるRNA配列を探したいとき、本発明によれば、以下の手順でこれを調べることができる。例として、入力構造トポロジーを T_2 、検査対象配列集合を s_1 、 s_2 とした場合について示す。

手順1) トポロジー(生成文法)を文法DBから選択する。ここでは T_2 (G_2)を選択する。

手順 2) 適合度の閾値を設定する。ここでは 5 を選択する。

手順 3) 検査対象とする RNA 配列集合を配列 DB から選択する。あるいは直接入力する。ここでは s_1 、 s_2 を選択する。

5 手順 4) 手順 3 で得られた各配列を、手順 1 で得られた生成文法によってそれぞれ構文解析を行ない、適合度最大の構文解析木をそれぞれ求める。ここでは s_1 を G_2 によって構文解析し、導出不能であることを得る。さらに s_2 を G_2 によって構文解析し、最大の適合度 6.7 を持つ構文解析木を得る。(第 17 図 (2) 参照)

10 手順 5) 手順 4 で得られた構文解析木のうち手順 2 で得られた閾値以上の適合度をもつ構文解析木に対応する配列を出力する。手順 4 で得られた s_2 の G_2 に対する適合度 6.7 の構文解析木は手順 2 で設定された閾値 5 よりも大きいので s_2 を出力する。以上の結果から、選択されたトポロジーをとりうる配列の候補が第 19 図のように出力される。

15 本発明の実施について厳密に上記の手順と同じである必要はない。例えば、手順 1 と 2 と 3 は任意の順にいれかえてもよいし、手順 5 は手順 4 の構文解析部分に含めてしまってもよい。

(3) 共通構造抽出

ある RNA 配列の集合に対して、これらの配列が共通してとりうる構造トポロジーを調べたいとき、本発明によれば、以下の手順でこれを調べることができる。
20 例として、入力配列集合を s_1 、 s_2 とし、検査対象構造トポロジーの集合を T_1 、 T_2 とした場合について示す。

手順 1) RNA 配列の集合を配列 DB から指定する。あるいは直接入力する。ここでは s_1 と s_2 を指定する。

25 手順 2) 検査対象とするトポロジーの集合 (生成文法の集合) を文法 DB から選択する。ここでは $T_1 (G_1)$ と $T_2 (G_2)$ を選択する。

手順 3) 適合度の閾値を設定する。閾値は手順 2 で得られた各トポロジー (生成文法) に対してそれぞれ設定してもよいし、共通にひとつ設定してもよい。こ

ここでは共通に0を設定する。

手順4) 手順1で得られた各配列を、手順2で得られた各生成文法によってそれぞれ構文解析を行ない、適合度最大の構文解析木を求める。

s_1 を G_1 によって構文解析し、最大の適合度1.4を持つ構文解析木を得る（

5 第14図（2）参照）。

s_1 を G_2 によって構文解析し、導出不能であることを得る。

s_2 を G_1 によって構文解析し、最大の適合度13.6を持つ構文解析木を得る。

（第17図（1）参照）

s_2 を G_2 によって構文解析し、最大の適合度6.7を持つ構文解析木を得る。

10 （第17図（2）参照）

手順5) 手順4で得られた構文解析木のうち閾値以上の適合度を持つ構文解析木を抽出する。手順4で得られたすべての構文解析木は手順3で得られた閾値0よりも大きい適合度を持つので手順4で得られたすべての構文解析木を抽出する。

手順6) 手順1で得られた配列集合を行に、手順2で得られたトポロジー集合
15 を列に、手順5で得られた構文解析木の適合度を要素に持つマトリックスを作成する。第20図に示すマトリックスを得る。

以上の結果得られたマトリックスを出力すれば、対象配列集合が共通してとりうる構造トポロジーを容易に確認することが可能になる。あるいは、以下の追加手順を実行すれば、共通構造の候補を順位づけて出力することができる。

20 手順7) 手順6で得られたマトリックスの各列、すなわちトポロジー、に対してスコアを計算する。例えば、有効な行要素の数を各列ごとに計算しスコアとすると、 T_1 のスコアは2、 T_2 のスコアは1となる。例えば、各行の適合度の総和を各列ごとに計算しスコアとすると、 T_1 のスコアは15.0、 T_2 のスコアは6.7となる。

25 手順8) 手順7で得られたスコアの高い順にトポロジーをソートし、出力する。上記のいずれのスコアを用いても T_1 、 T_2 の順に出力される。

また、本発明の実施について厳密に上記の手順と同じである必要はない。例え

ば、手順 1 と 2 は順序をいれかえてもよいし、手順 5 を手順 4 の構文解析部分に含めてしまってもよい。

(4) ジーン・ファインダ

RNA 遺伝子部分に対応する配列は、非常に安定な構造をとりやすいので、適合度が高くなる。そこで本発明では、万能文法を用いて構文解析を行ない、適合度の高い配列を配列 DB から選び出して遺伝子候補として出力する。例として、配列集合を s_1 、 s_2 とした場合について示す。

手順 1) 検査対象とする RNA 配列の集合を配列 DB から指定する。あるいは直接入力する。ここでは s_1 と s_2 を指定する。

10 手順 2) 適合度の閾値を設定する。ここでは 10 を設定する。

手順 3) 手順 1 で得られた各配列を万能文法 G_0 によってそれぞれ構文解析を行ない、適合度最大の構文解析木を求める。

s_1 を G_0 によって構文解析し、最大の適合度 1.4 を持つ構文解析木を得る。

s_2 を G_0 によって構文解析し、最大の適合度 13.6 を持つ構文解析木を得る。

15 手順 4) 手順 3 で得られた構文解析木のうち閾値以上の適合度をもつ構文解析木に対応する配列を遺伝子候補として出力する。手順 3 で得られた s_1 の構文解析木は閾値 10 に満たないので s_1 は出力しない。手順 3 で得られた s_2 の構文解析木は閾値 10 よりも大きいので s_2 を遺伝子候補として出力する。

20 本発明の実施について厳密に上記の手順と同じである必要はない。例えば、手順 1 と 2 は順序をいれかえてもよいし、手順 4 は手順 3 の構文解析部分に含めてしまってもよい。

(5) RNA 配列から同じ構造を持つ RNA 配列を出力

ある RNA 配列集合に対して、これらと同じトポロジーをとりうる RNA 配列を調べたいとき、(3) の発明と (2) の発明とを組み合わせた本発明によれば、
25 以下の手順でこれを調べることができる。例として、入力配列を $s = g c c c a a a g g g c a g c c c a a a g g g c$ 、検査対象トポロジー集合を T_1 、 T_2 、検査対象配列集合を s_1 、 s_2 とした場合について示す。

手順 1) RNA 配列集合を入力する。ここでは s のみからなる配列集合を入力する。

手順 2) 検査対象とする RNA 配列の集合を配列 DB から指定する。ここでは s_1 と s_2 を指定する。

- 5 手順 3) 検査対象とするトポロジーの集合（生成文法の集合）を文法 DB から選択する。ここでは $T_1 (G_1)$ と $T_2 (G_2)$ を選択する。

手順 4) 適合度の閾値を設定する。閾値は手順 3 で得られた各トポロジー（生成文法）に対してそれぞれ設定してもよいし、共通にひとつ設定してもよい。ここでは共通に 5 を設定する。

- 10 手順 5) 手順 1 で得られた各 RNA 配列を、手順 2 で得られた各生成文法によってそれぞれ構文解析を行ない、適合度最大の構文解析木をそれぞれ求める。ここでは s を G_1 によって構文解析し、最大の適合度 3. 1 を持つ構文解析木を得る。第 21 図 (1) にこの構文解析木が表現する二次構造を示す。さらに s を G_2 によって構文解析し、最大の適合度 5. 1 を持つ構文解析木を得る。第 21 図
15 (2) にこの構文解析木が表現する二次構造を示す。

手順 6) 手順 5 で得られた構文解析木のうち、手順 4 で得られた閾値以上の適合度をもつ構文解析木に対応する構文解析木を抽出する。手順 5 で得られた構文解析木のうち、 G_2 で構文解析することによって得られた適合度 5. 1 の構文解析木が手順 4 で得られた閾値 5 よりも大きいのでこれを抽出する。

- 20 手順 7) 手順 1 で得られた配列集合を行に、手順 3 で得られたトポロジー集合を列に、手順 6 で得られた構文解析木の適合度を要素に持つマトリックスを作成する。第 22 図に示すマトリックスを得る。

- 手順 8) 手順 6 で得られたマトリックスの各列、すなわちトポロジー、に対してスコアを計算し、スコアの順にトポロジーをソートする。ここでは行の総和を
25 各列ごとに計算しスコアとするが、1 行しかないので結果として、 T_1 のスコアが未定義、 T_2 のスコアが 5. 1 になる。スコアを持つものだけでソートすると、 T_2 のみ得られる。

手順 9) 手順 2 で得られた各配列を、手順 8 で得られたトポロジーの順にそれぞれ対応する文法で構文解析を行ない、適合度最大の構文解析木をそれぞれ求める。ここでは s_1 を G_2 によって構文解析し、導出不能であることを得る。

さらに s_2 を G_2 によって構文解析し、最大の適合度 6. 7 を持つ構文解析木を得る。(第 17 図 (2) 参照)

手順 10) 手順 9 で得られた構文解析木のうち手順 4 で得られた閾値以上の適合度をもつ構文解析木に対応する配列を出力する。このとき、あわせてトポロジーとそのトポロジーに対する手順 8 で得られたスコアを出力する。手順 9 で得られた s_2 の G_2 に対する構文解析木の適合度 6. 7 は手順 4 で得られた閾値 5 よりも大きいので s_2 を出力する。あわせて、 T_2 とそのスコア 5. 1 を出力する。

以上の結果から第 23 図に示すような出力が得られる。

この結果、 s_2 がトポロジー T_2 において、 s と共通な構造をとりうるようになる。

本発明の実施について厳密に上記の手順と同じである必要はない。例えば、手順 1 と 2 と 3 は任意の順に入れ換えてもよいし、手順 6 は手順 5 の構文解析部分に含めてしまってもよいし、手順 10 の閾値により構文解析木を取捨選択する部分は手順 9 の構文解析部分に含めてしまってもよい。

[他の実施の形態]

さて、これまで本発明の実施の形態について説明したが、本発明は、上述した実施の形態以外にも、上記特許請求の範囲に記載した技術的思想の範囲内において種々の異なる実施の形態にて実施されてよいものである。

例えば、RNA 配列解析装置 100 がスタンドアローンの形態で RNA 配列解析方法を行う場合を一例に説明したが、RNA 配列解析装置 100 とは別筐体で構成されるクライアント端末からの要求に応じて RNA 配列解析方法を行い、その処理結果を当該クライアント端末に返却するように構成してもよい。

また、構造予測部 102 a は、適合度計算部 102 c による適合度計算を行いながら構文解析部 102 b により構文解析木を導出してもよい。すなわち、構文

解析木を導出する構文解析部 102b と、導出された構文解析木に対して適合度の計算を行う適合度計算部 102c をひとつのアルゴリズムにて実現してもよい。このように構成することにより、RNA 配列と木文法に対して可能な構文解析木は無数（配列長に対して指数のオーダー）に存在するため、構文解析木を導出し

5 てから適合度計算を行いソートすると指数オーダーの計算時間と記憶容量が必要となるという問題点を解決することができる。

また、実施の形態において説明した各処理のうち、自動的に行なわれるものとして説明した処理の全部または一部を手動的に行うこともでき、あるいは、手動的に行なわれるものとして説明した処理の全部または一部を公知の方法で自動的

10 に行うこともできる。

特に、構造予測部 102a は複数のタスクとして実現してもよく、それぞれのタスクで並列処理を行うように実現してもよい。

この他、上記文書中や図面中で示した処理手順、制御手順、具体的名称、各種の登録データや検索条件等のパラメータを含む情報、画面例、データベース構成

15 については、特記する場合を除いて任意に変更することができる。

また、RNA 配列解析装置 100 に関して、図示の各構成要素は機能概念的なものであり、必ずしも物理的に図示の如く構成されていることを要しない。

例えば、RNA 配列解析装置 100 の各サーバが備える処理機能、特に制御部にて行なわれる各処理機能については、その全部または任意の一部を、CPU（

20 C e n t r a l P r o c e s s i n g U n i t ）および当該 CPU にて解釈実行されるプログラムにて実現することができ、あるいは、ワイヤードロジックによるハードウェアとして実現することも可能である。なお、プログラムは、後述する記録媒体に記録されており、必要に応じて RNA 配列解析装置 100 に機械的に読み取られる。

25 記憶部 106 に格納される各種のデータベース等（RNA 配列データベース 106a ～共通構造マトリックス 106c）は、RAM、ROM 等のメモリ装置、ハードディスク等の固定ディスク装置、フレキシブルディスク、光ディスク等の

ストレージ手段であり、各種処理やウェブサイト提供に用いる各種のプログラムやテーブルやファイルやデータベースやウェブページ用ファイル等を格納する。

また、RNA配列解析装置100は、既知のパーソナルコンピュータ、ワークステーション等の情報処理端末等の情報処理装置にプリンタやモニタやイメージ
5 スキャナ等の周辺装置を接続し、該情報処理装置に本発明の方法を実現させるソフトウェア（プログラム、データ等を含む）を実装することにより実現してもよい。

さらに、RNA配列解析装置100の分散・統合の具体的形態は図示のものに限られず、その全部または一部を、各種の負荷等に応じた任意の単位で、機能的
10 または物理的に分散・統合して構成することができる。例えば、各データベースを独立したデータベース装置として独立に構成してもよく、また、処理の一部をCGI（Common Gateway Interface）を用いて実現してもよい。

また、本発明にかかるプログラムを、コンピュータ読み取り可能な記録媒体に
15 格納することもできる。ここで、この「記録媒体」とは、フレキシブルディスク、光磁気ディスク、ROM、EPROM、EEPROM、CD-ROM、MO、DVD等の任意の「可搬用の物理媒体」や、各種コンピュータシステムに内蔵されるROM、RAM、HD等の任意の「固定用の物理媒体」、あるいは、LAN、WAN、インターネットに代表されるネットワークを介してプログラムを送信する
20 場合の通信回線や搬送波のように、短期にプログラムを保持する「通信媒体」を含むものとする。

また、「プログラム」とは、任意の言語や記述方法にて記述されたデータ処理方法であり、ソースコードやバイナリコード等の形式を問わない。なお、「プログラム」は必ずしも単一的に構成されるものに限られず、複数のモジュールやライブラリとして分散構成されるものや、OS（Operating System）に代表される別個のプログラムと協働してその機能を達成するものをも含む。
25 なお、実施の形態に示した各装置において記録媒体を読み取るための具体的な構

成、読み取り手順、あるいは、読み取り後のインストール手順等については、周知の構成や手順を用いることができる。

また、ネットワーク 300 は、RNA 配列解析装置 100 と外部システム 200 とを相互に接続する機能を有し、例えば、インターネットや、イントラネット
5 や、LAN（有線／無線の双方を含む）や、VAN や、パソコン通信網や、公衆電話網（アナログ／デジタルの双方を含む）や、専用回線網（アナログ／デジタルの双方を含む）や、CATV 網や、IMT 2000 方式、GSM 方式または PDC／PDC-P 方式等の携帯回線交換網／携帯パケット交換網や、無線呼出網
10 や、Bluetooth 等の局所無線網や、PHS 網や、CS、BS または ISDB 等の衛星通信網等のうちいずれかを含んでもよい。すなわち、本システムは、有線・無線を問わず任意のネットワークを介して、各種データを送受信することができる。

以上詳細に説明したように、本発明によれば、RNA 二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA 配
15 列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度を計算し、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を適合度が高い順にソートし、ソートされた構文解析木を RNA 配列の二次構造の候補として出力するので、一配列に対して多文法で構文解析を行うことができるようになる。すなわち、各生成文法に対してそれぞれ構文解析し適合
20 度計算を行い適合度を得る。その結果、生成文法ごとに適合度が得られることになり、それらの適合度をソートすることによって生成文法に順位を付ける。これにより、生成文法に対する構造トポロジーにも順位が付けられることになるので、最終的に RNA 配列が取り得る可能性の高い順に構造トポロジーを確認することができる RNA 配列解析装置、RNA 配列解析方法、プログラム、および、記録
25 媒体を提供することができる。

また、本発明によれば、RNA 二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA 配列を生成文法に適用し

て構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力するので、多配列に対して一文法で構文解析を行うことができるようになる。

- 5 すなわち、与えられた特定の構造トポロジーに対し、対応する生成文法を取得し、これを用いてRNA配列データベースに格納されているすべてまたは一部のRNA配列をそれぞれ構文解析し、ある閾値以下の適合度で構文解析に成功したRNA配列群を結果として出力する。これにより、与えられた特定の二次構造を取り得るようなRNA配列を検索することができるRNA配列解析装置、RNA配列解析方法、プログラム、および、記録媒体を提供することができる。
- 10

- また、本発明によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を抽出し、構造トポロジーとRNA配列とを二次元マトリックスで表示し、二次元マトリックスにおいて抽出されたRNA配列と構造トポロジーに対応する格子部分をマークすることにより、RNA配列間で共通に有する構造トポロジーを可視化するので、RNA配列間の共通構造を容易に発見することができるRNA配列解析装置、RNA配列解析方法、プログラム、および、記録媒体を提供することができる。
- 15
- 20

- また、本発明によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、利用者が入力したDNA配列から転写されるRNA配列を作成し、作成されたRNA配列に対して生成文法を適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列に対応するDNA配列部分を遺伝子の候補として予測するので、DNA配列のうち既知のトポロジーを有する可能性のあるRNA配列に対応する
- 25

部分について、遺伝子部分である可能性があることを予測することができるRNA配列解析装置、RNA配列解析方法、プログラム、および、記録媒体を提供することができる。

また、本発明によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度に基づいてRNA配列間の類似度を計算するので、RNA構造の類似度を容易に求めることができるRNA配列解析装置、RNA配列解析方法、プログラム、および、記録媒体を提供することができる。

さらに、本発明によれば、RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納し、RNA配列を生成文法に適用して構文解析木を導出し、導出された構文解析木に対して適合度の計算を行い、計算された適合度のうち予め定めた条件を満たす適合度である構文解析木を導出したRNA配列を抽出し、構造トポロジーとRNA配列とを二次元マトリックスで表示し、二次元マトリックスにおいて抽出されたRNA配列と構造トポロジーに対応する格子部分に適合度を表示する適合度マトリックスを作成し、適合度マトリックスについて、適合度により構造トポロジーをソートし、他のRNA配列について当該ソートされた構造トポロジーの順番に対応する生成文法により構文解析を行い適合度が最大となる構文解析木を求め、予め定めた条件を満たす適合度を持つ構文解析木に対応する他のRNA配列を抽出するので、共通の構造を持つRNA配列を容易に発見することができるRNA配列解析装置、RNA配列解析方法、プログラム、および、記録媒体を提供することができる。

産業上の利用可能性

以上のように、本発明にかかるRNA配列解析装置、RNA配列解析方法、プログラム、および、記録媒体は、RNA二次構造予測、RNA配列解析、遺伝子予測、および、これらを利用した創薬等に適している。

請 求 の 範 囲

1. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、
- 5 RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段と、
上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、
上記適合度計算手段により計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を上記適合度が高い順にソートするソート手段と、
- 10 上記ソート手段によりソートされた上記構文解析木を上記RNA配列の二次構造の候補として出力する出力手段と、
を備えたことを特徴とするRNA配列解析装置。
2. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成
- 15 文法とを対応付けて格納する文法格納手段と、
RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段と、
上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、
上記適合度計算手段により計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を上記構造トポロジー
- 20 の二次構造を持つRNA配列の候補として出力する出力手段と、
を備えたことを特徴とするRNA配列解析装置。
3. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成
- 25 文法とを対応付けて格納する文法格納手段と、
RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段と、
上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う

適合度計算手段と、

上記適合度計算手段により計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出手段と、

- 上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記
- 5 二次元マトリックスにおいて上記抽出手段にて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分をマークすることにより、上記RNA配列間で共通に有する構造トポロジーを可視化する共通構造マトリックス作成手段と、
- を備えたことを特徴とするRNA配列解析装置。

- 10 4. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、

利用者が入力したDNA配列から転写されるRNA配列を作成するRNA配列作成手段と、

- 上記RNA配列作成手段により作成された上記RNA配列に対して上記生成文
- 15 法を適用して構文解析木を導出する構文解析手段と、

上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、

- 上記適合度計算手段により計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列に対応する上記DNA
- 20 配列部分を遺伝子の候補として予測する遺伝子予測手段と、
- を備えたことを特徴とするRNA配列解析装置。

5. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、

- 25 RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段と、
- 上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、

上記適合度計算手段により計算された上記適合度に基づいて上記RNA配列間の類似度を計算する類似度計算手段と、

を備えたことを特徴とするRNA配列解析装置。

- 5 6. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納手段と、

RNA配列を上記生成文法に適用して構文解析木を導出する構文解析手段と、

上記構文解析手段にて導出された上記構文解析木に対して適合度の計算を行う適合度計算手段と、

- 10 上記適合度計算手段により計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出手段と、

上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出手段にて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分に上記適合度を表示する適合度マトリックス

- 15 を作成する適合度マトリックス作成手段と、

上記適合度マトリックス作成手段にて作成された上記適合度マトリックスについて、上記適合度により上記構造トポロジーをソートし、他のRNA配列について当該ソートされた上記構造トポロジーの順番に対応する上記生成文法により構文解析を行い上記適合度が最大となる上記構文解析木を求め、予め定めた条件を満たす上記適合度を持つ上記構文解析木に対応する上記他のRNA配列を抽出する共通構造抽出手段と、

- 20 を備えたことを特徴とするRNA配列解析装置。

7. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

- 25 RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を上記適合度が高い順にソートするソートステップと、

上記ソートステップによりソートされた上記構文解析木を上記RNA配列の二次構造の候補として出力する出力ステップと、
を含むことを特徴とするRNA配列解析方法。

10 8. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

15 上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力する出力ステップと、
を含むことを特徴とするRNA配列解析方法。

20

9. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

25 上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を

満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出ステップと、

- 上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出ステップにて抽出された上記RNA配列と
- 5 上記構造トポロジーに対応する格子部分をマークすることにより、上記RNA配列間で共通に有する構造トポロジーを可視化する共通構造マトリックス作成ステップと、
- を含むことを特徴とするRNA配列解析方法。

- 10 10. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

利用者が入力したDNA配列から転写されるRNA配列を作成するRNA配列作成ステップと、

- 上記RNA配列作成ステップにより作成された上記RNA配列に対して上記生
- 15 成文法を適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

- 上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列に対応する上記D
- 20 NA配列部分を遺伝子の候補として予測する遺伝子予測ステップと、
- を含むことを特徴とするRNA配列解析方法。

11. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

- 25 RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を

行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度に基づいて上記RNA配列間の類似度を計算する類似度計算ステップと、
を含むことを特徴とするRNA配列解析方法。

5

12. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

10 上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出ステップと、

15 上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出ステップにて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分に上記適合度を表示する適合度マトリックスを作成する適合度マトリックス作成ステップと、

上記適合度マトリックス作成ステップにて作成された上記適合度マトリックス
20 について、上記適合度により上記構造トポロジーをソートし、他のRNA配列について当該ソートされた上記構造トポロジーの順番に対応する上記生成文法により構文解析を行い上記適合度が最大となる上記構文解析木を求め、予め定めた条件を満たす上記適合度を持つ上記構文解析木に対応する上記他のRNA配列を抽出する共通構造抽出ステップと、

25 を含むことを特徴とするRNA配列解析方法。

13. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生

成文法とを対応付けて格納する文法格納ステップと、

RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を上記適合度が高い順にソートするソートステップと、

上記ソートステップによりソートされた上記構文解析木を上記RNA配列の二次構造の候補として出力する出力ステップと、

を含むことを特徴とするRNA配列解析方法をコンピュータに実行させるプログラム。

14. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力する出力ステップと、

を含むことを特徴とするRNA配列解析方法をコンピュータに実行させるプログラム。

25

15. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

- 5 上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出ステップと、

上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出ステップにて抽出された上記RNA配列と
10 上記構造トポロジーに対応する格子部分をマークすることにより、上記RNA配列間で共通に有する構造トポロジーを可視化する共通構造マトリックス作成ステップと、

を含むことを特徴とするRNA配列解析方法をコンピュータに実行させるプログラム。

15

16. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

利用者が入力したDNA配列から転写されるRNA配列を作成するRNA配列作成ステップと、

- 20 上記RNA配列作成ステップにより作成された上記RNA配列に対して上記生成文法を適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

- 25 上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列に対応する上記DNA配列部分を遺伝子の候補として予測する遺伝子予測ステップと、

を含むことを特徴とするRNA配列解析方法をコンピュータに実行させるプロ

グラム。

17. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

5 RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

10 上記適合度計算ステップにより計算された上記適合度に基づいて上記RNA配列間の類似度を計算する類似度計算ステップと、

を含むことを特徴とするRNA配列解析方法をコンピュータに実行させるプログラム。

18. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

15 RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

20 上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出ステップと、

上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出ステップにて抽出された上記RNA配列と
25 上記構造トポロジーに対応する格子部分に上記適合度を表示する適合度マトリックスを作成する適合度マトリックス作成ステップと、

上記適合度マトリックス作成ステップにて作成された上記適合度マトリックス

について、上記適合度により上記構造トポロジーをソートし、他のRNA配列について当該ソートされた上記構造トポロジーの順番に対応する上記生成文法により構文解析を行い上記適合度が最大となる上記構文解析木を求め、予め定めた条件を満たす上記適合度を持つ上記構文解析木に対応する上記他のRNA配列を抽出

5 出する共通構造抽出ステップと、

を含むことを特徴とするRNA配列解析方法をコンピュータに実行させるプログラム。

19. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

15 上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を上記適合度が高い順にソートするソートステップと、

上記ソートステップによりソートされた上記構文解析木を上記RNA配列の二次構造の候補として出力する出力ステップと、

20 を含むRNA配列解析方法をコンピュータに実行させるプログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

20. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

25 RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を

行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を上記構造トポロジーの二次構造を持つRNA配列の候補として出力する出力ステップと、

- 5 を含むRNA配列解析方法をコンピュータに実行させるプログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

21. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

- 10 RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

- 15 上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出ステップと、

- 20 上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出ステップにて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分をマークすることにより、上記RNA配列間で共通に有する構造トポロジーを可視化する共通構造マトリックス作成ステップと、

を含むRNA配列解析方法をコンピュータに実行させるプログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

- 25 22. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

利用者が入力したDNA配列から転写されるRNA配列を作成するRNA配列

作成ステップと、

上記RNA配列作成ステップにより作成された上記RNA配列に対して上記生成文法を適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列に対応する上記DNA配列部分を遺伝子の候補として予測する遺伝子予測ステップと、

を含むRNA配列解析方法をコンピュータに実行させるプログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

23. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度に基づいて上記RNA配列間の類似度を計算する類似度計算ステップと、

を含むRNA配列解析方法をコンピュータに実行させるプログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

24. RNA二次構造の構造トポロジーと、当該構造トポロジーに適合する生成文法とを対応付けて格納する文法格納ステップと、

RNA配列を上記生成文法に適用して構文解析木を導出する構文解析ステップと、

上記構文解析ステップにて導出された上記構文解析木に対して適合度の計算を

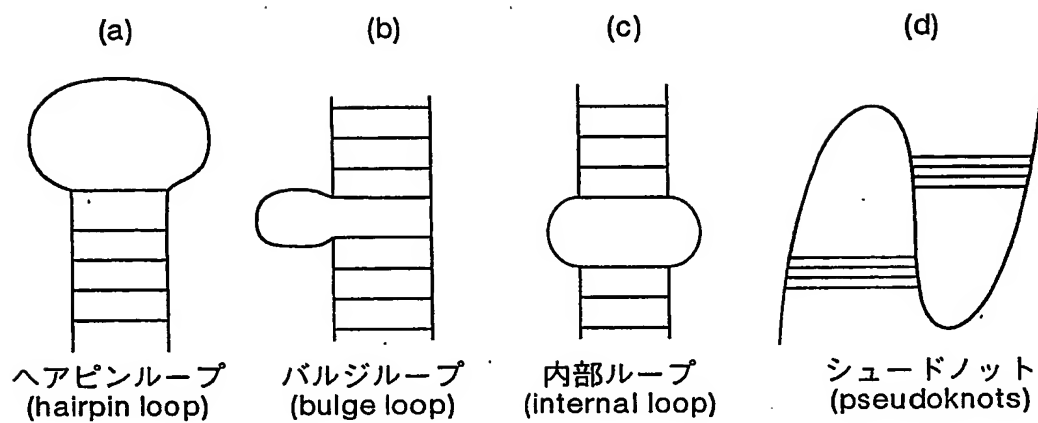
行う適合度計算ステップと、

上記適合度計算ステップにより計算された上記適合度のうち予め定めた条件を満たす適合度である上記構文解析木を導出した上記RNA配列を抽出する抽出ステップと、

- 5 上記構造トポロジーと上記RNA配列とを二次元マトリックスで表示し、上記二次元マトリックスにおいて上記抽出ステップにて抽出された上記RNA配列と上記構造トポロジーに対応する格子部分に上記適合度を表示する適合度マトリックスを作成する適合度マトリックス作成ステップと、

- 10 上記適合度マトリックス作成ステップにて作成された上記適合度マトリックスについて、上記適合度により上記構造トポロジーをソートし、他のRNA配列について当該ソートされた上記構造トポロジーの順番に対応する上記生成文法により構文解析を行い上記適合度が最大となる上記構文解析木を求め、予め定めた条件を満たす上記適合度を持つ上記構文解析木に対応する上記他のRNA配列を抽出する共通構造抽出ステップと、

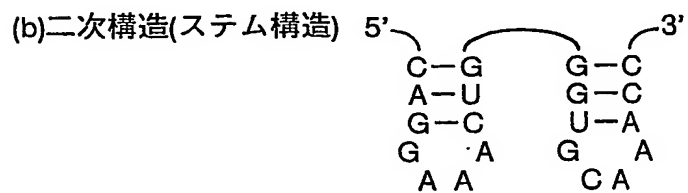
- 15 を含むRNA配列解析方法をコンピュータに実行させるプログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

1 / 2 3
第 1 図

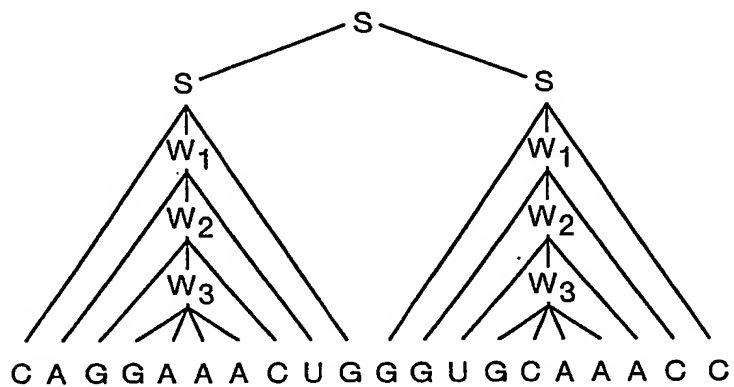
2 / 23
第 2 図

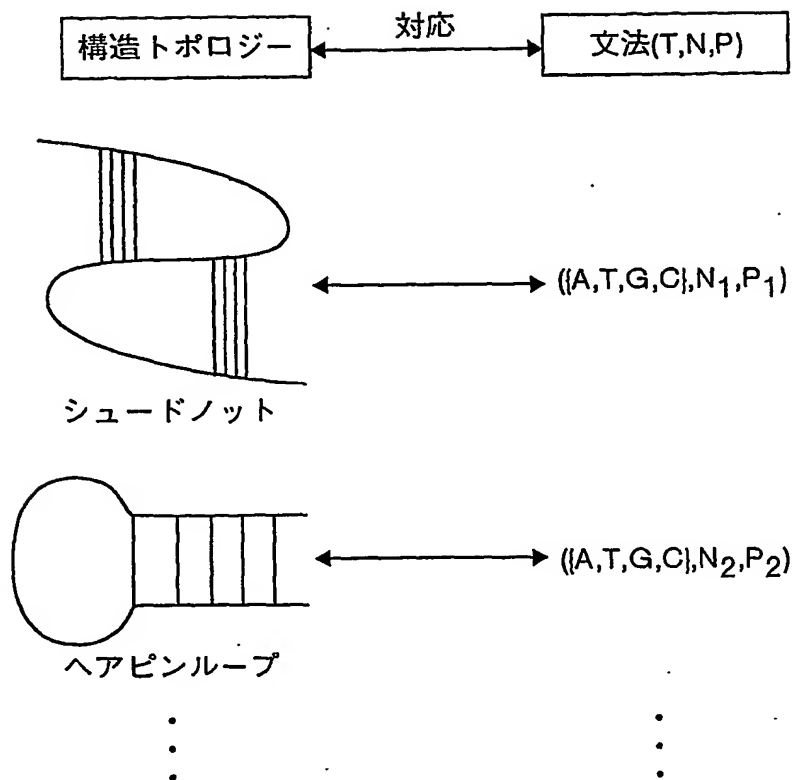
(a) RNA配列

C	A	G	G	A	A	A	C	U	G	G	G	U	G	C	A	A	A	C	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

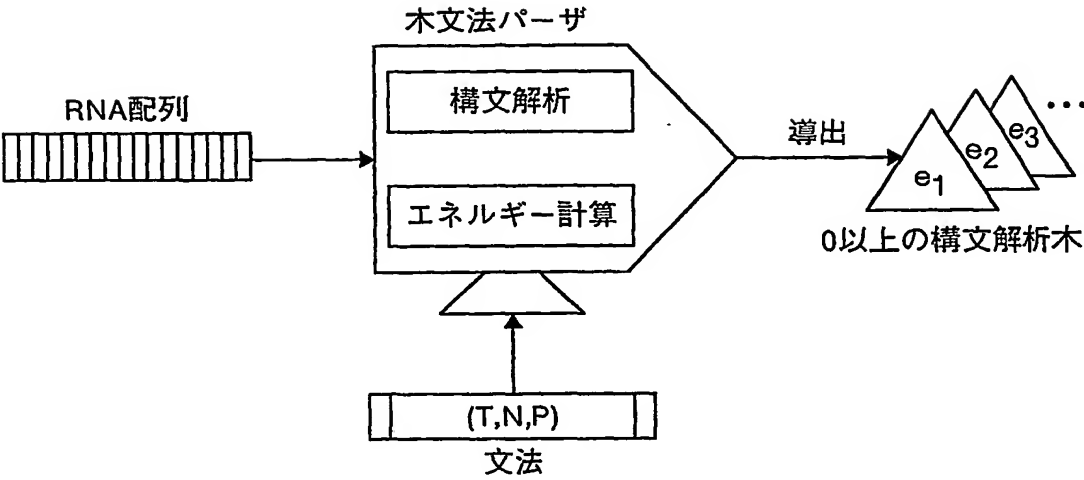


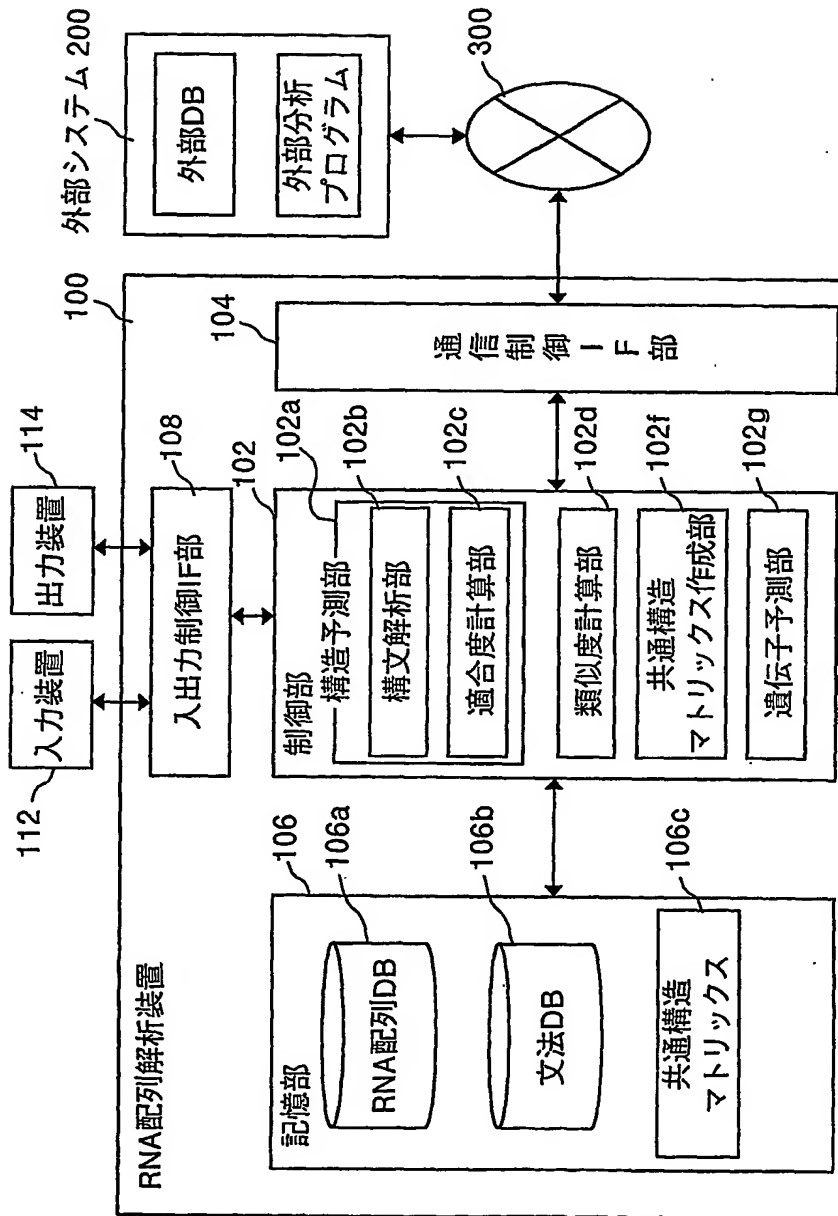
(c) 構文解析木



3 / 23
第3図

4 / 2 3
第 4 図

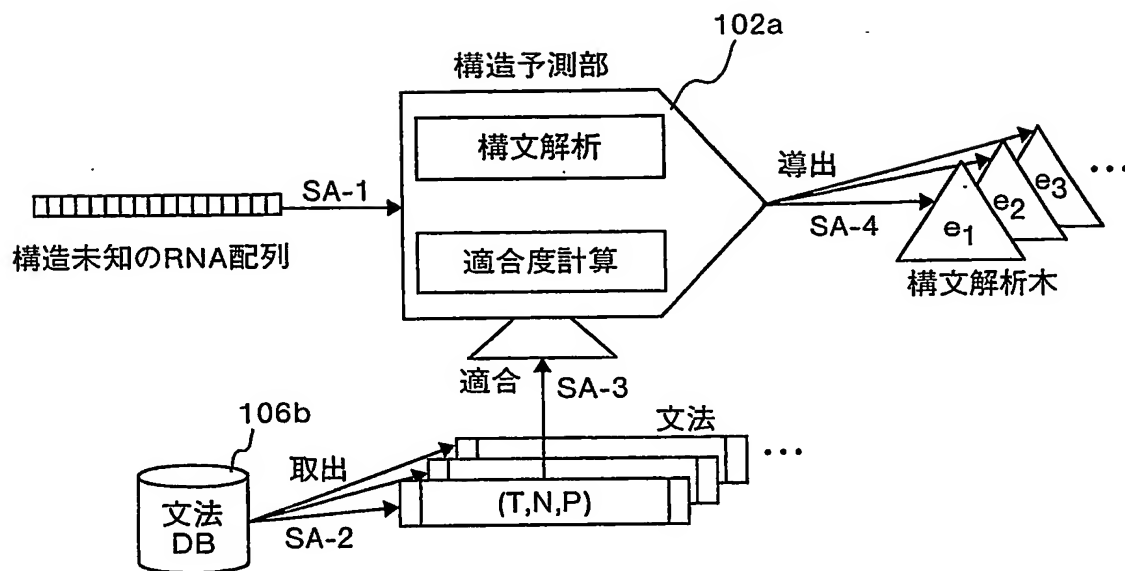


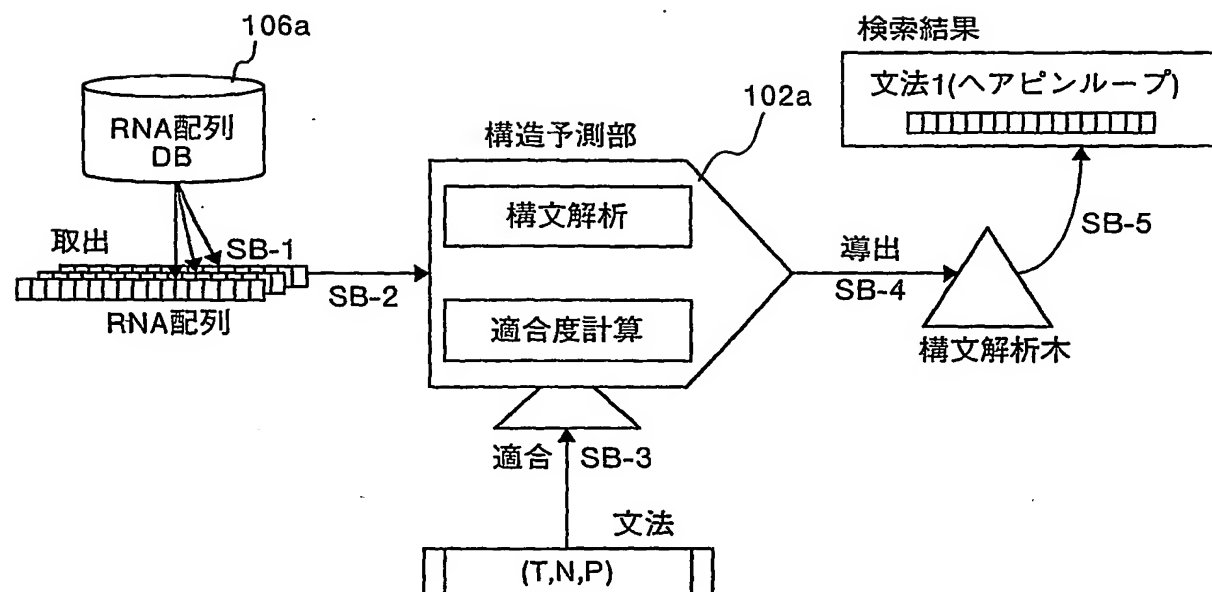
5 / 23
第5図

6 / 23
第6図

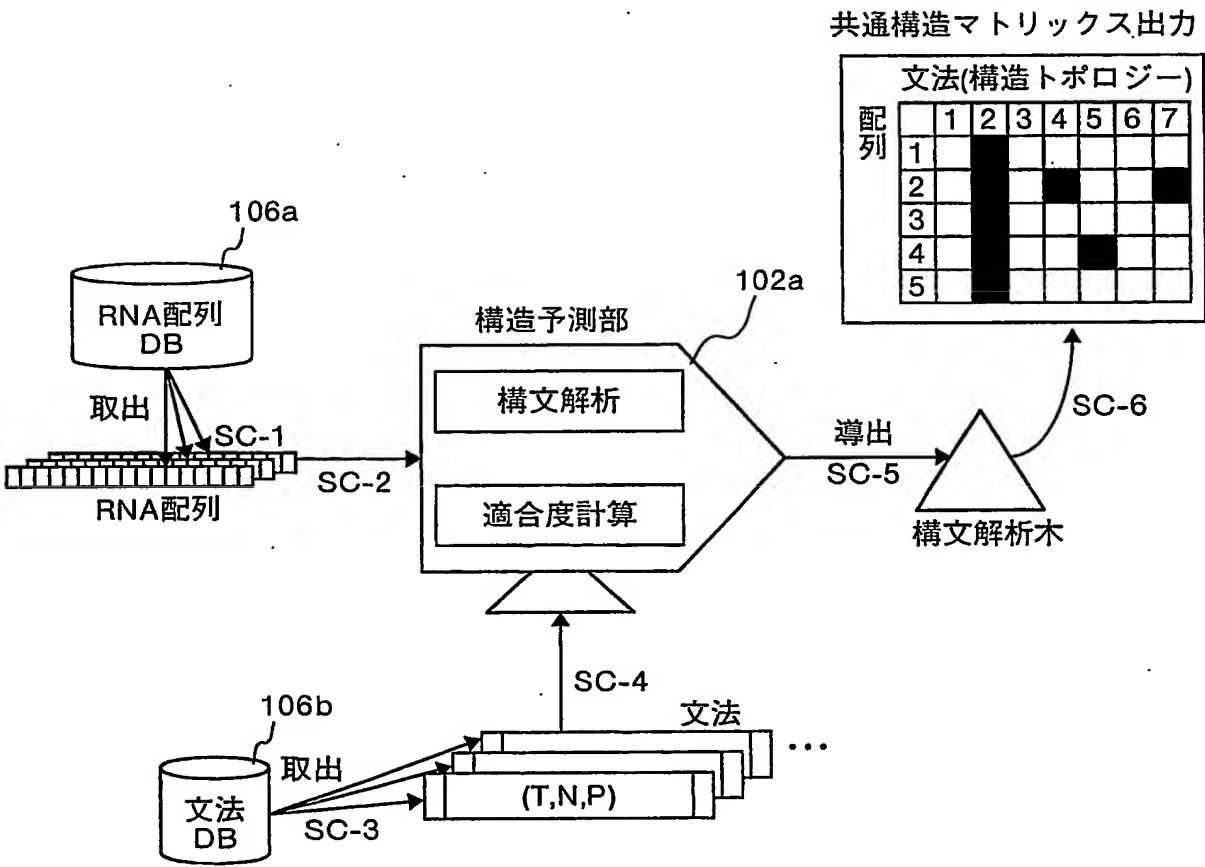
文法データベース 106b

構造トポロジ	文法		
	終端記号T	非終端記号N	生成規則P
シュードノット	{A,T,G,C}	N ₁	P ₁
ヘアピンループ		N ₂	P ₂
内側ループ		N ₃	P ₃
多枝ループ		N ₄	P ₄
バルジループ		N ₅	P ₅
⋮		⋮	⋮

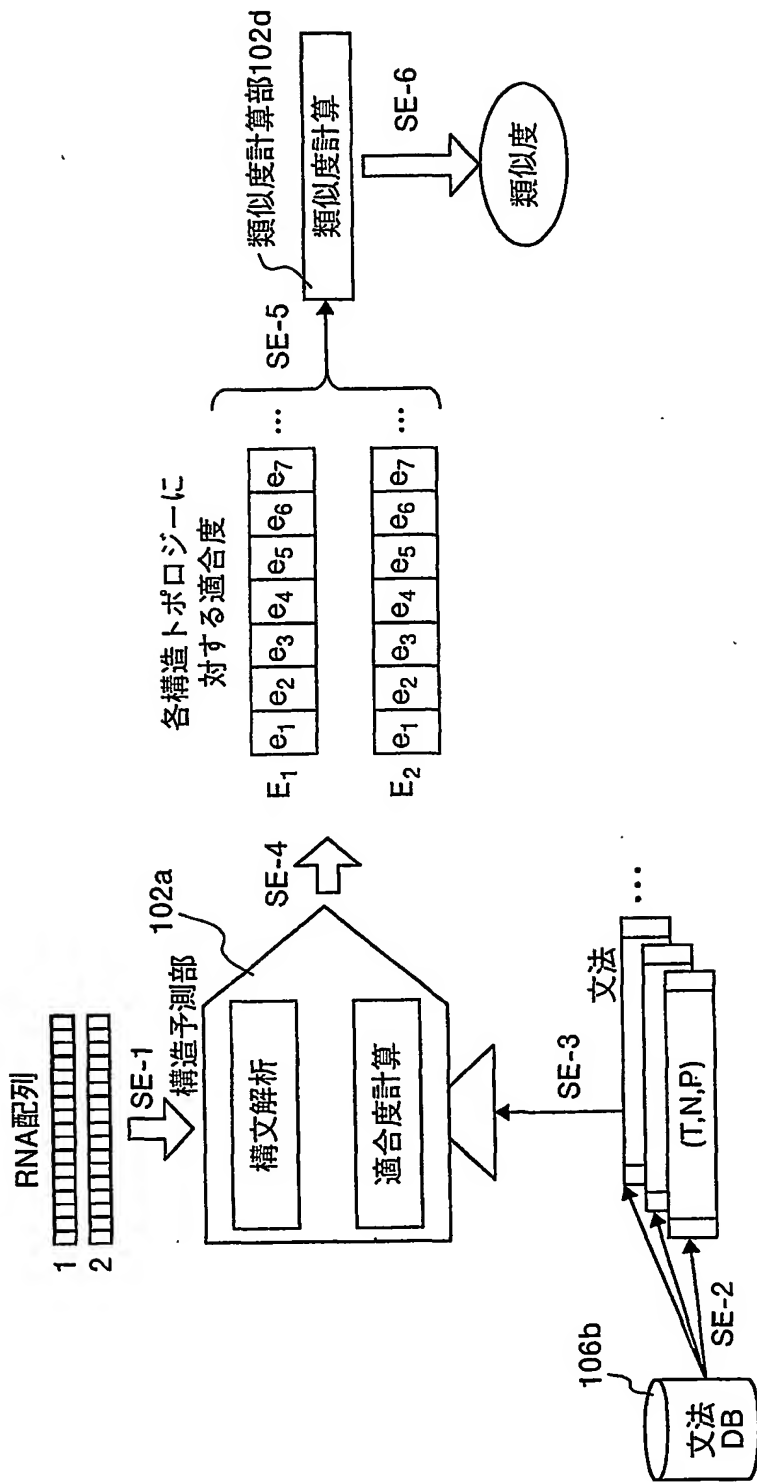
7 / 23
第7図

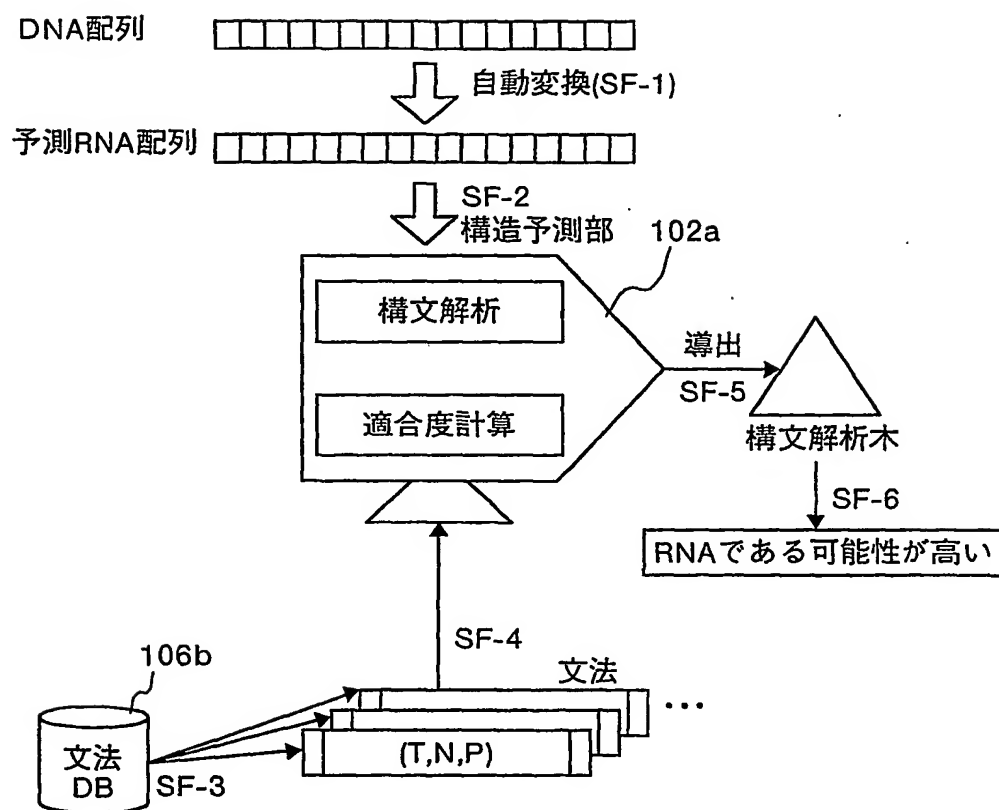
8/23
第8図

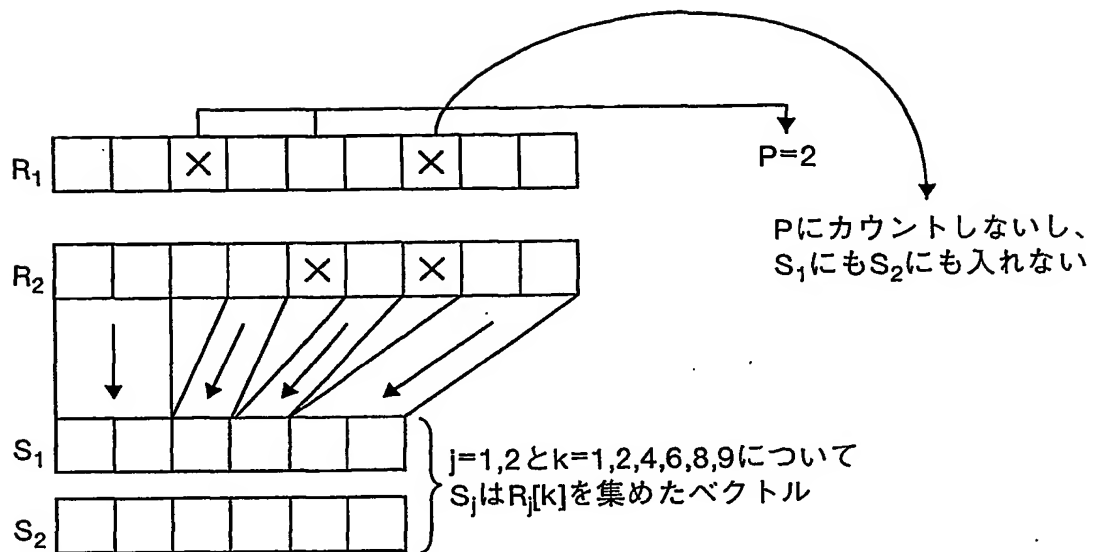
9 / 23
第9図



10/23
第10図

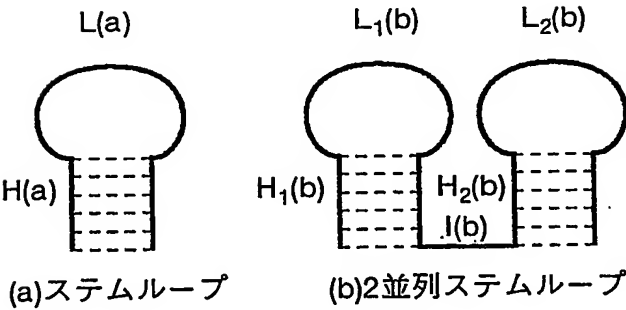


11/23
第11図

1 2 / 2 3
第 1 2 図

1 3 / 2 3
第 1 3 図

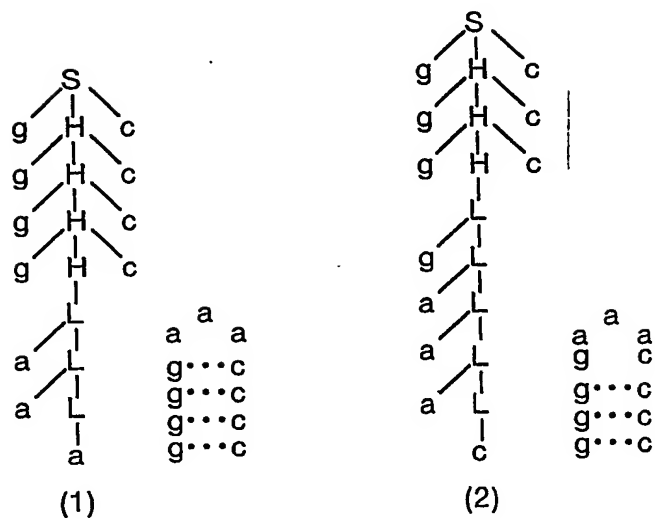
(RNA二次構造トポロジー例)



$$1 \ 4 \ / \ 2 \ 3$$

第 14 図

(s_1 の構文解析木と二次構造)



15/23
第15図

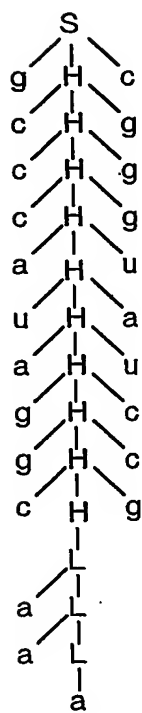
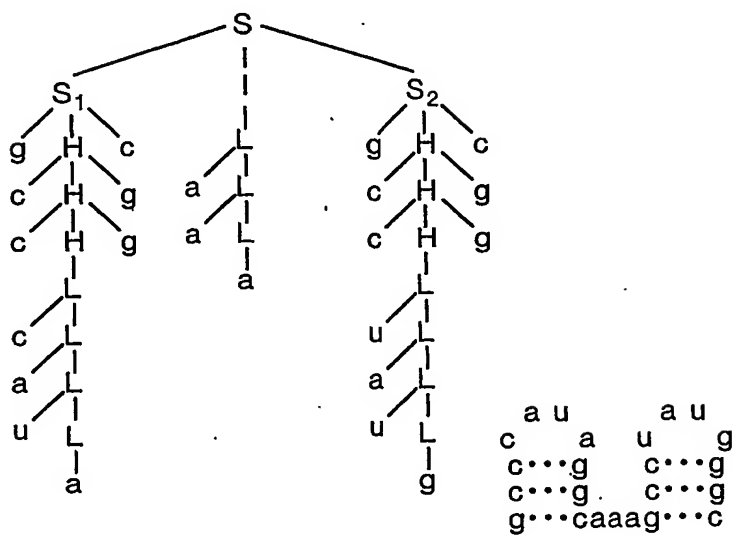
(塩基対の自由エネルギー[kcal/mol])

3'側塩基対	5'側塩基対				
	gu	au	ua	cg	gc
gu	-0.5	-0.5	-0.7	-1.5	-1.3
au	-0.5	-0.9	-1.1	-1.8	-2.3
ua	-0.7	-0.9	-0.9	-1.7	-2.1
cg	-1.9	-2.1	-2.3	-2.9	-3.4
gc	-1.5	-1.7	-1.8	-2.0	-2.9

(ループの自由エネルギー[kcal/mol])

ループ種別	ループサイズ													
	1	2	3	4	5	6	7	8	9	10	12	14	16	
バルジループ	3.3	5.2	6.0	6.7	7.4	8.2	9.1	10.0	10.5	11.0	11.8	12.5	13.0	
ヘアピンループ	-	-	7.4	5.9	4.4	4.3	4.1	4.1	4.2	4.3	4.9	5.6	6.1	
内部ループ	-	0.8	1.3	1.7	2.1	2.5	2.6	2.8	3.1	3.6	4.4	5.1	5.6	

17/23
第17図

(1)G₁による導出(s₂の最適な構文解析木と二次構造)(2)G₂による導出

18 / 23

第 18 図

順位	トポロジー	適合度	二次構造
1	ステムループ T_1	13.7	図18(1)に示された二次構造
2	2 並列ステムループ T_2	6.7	図18(2)に示された二次構造

19 / 23
第 19 図

配列	適合度	二次構造
s ₂	6.7	図18(2)に示された二次構造

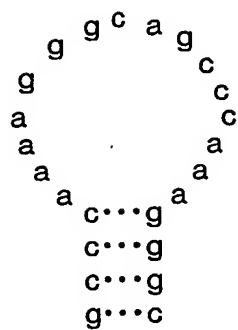
20/23

第20図

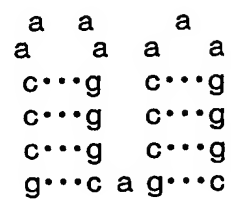
	T_1	T_2
s_1	1.4	-
s_2	13.6	6.7

21/23
第21図

(sの最適な二次構造)



(1)



(2)

22/23
第22図

	T_1	T_2
s	-	5.1

23 / 23
第 23 図

配列	適合度	トポロジー(スコア)
s ₂	6.7	T ₂ (5.1)

INTERNATIONAL SEARCH REPORT

International Publication No.

PCT/JP03/00011

A. CLASSIFICATION OF SUBJECT MATTER

Int.Cl⁷ G06F17/30, C12Q1/68, C12N15/00, G01N33/50

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

Int.Cl⁷ G06F17/30, C12Q1/68, C12N15/00, G01N33/50

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Jitsuyo Shinan Koho	1922-1996	Toroku Jitsuyo Shinan Koho	1994-2003
Kokai Jitsuyo Shinan Koho	1971-2003	Jitsuyo Shinan Toroku Koho	1996-2003

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

JICST FILE (JOIS), WPI, INSPEC (DIALOG)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	YOSHIZAWA et al., "RNA Niji Kozo Yosoku Shien System no Kaihatsu to Jisso", The Society for Artificial Intelligence Jinko Chino Kisoron Kenkyukai (Dai 32 Kai) Shiryo, 26 March, 1998 (26.03.98), pages 29 to 35, particularly, page 31	1, 2, 4, 5, 7, 8, 10, 11, 13, 14, 16, 17, 19, 20, 22, 23, 3, 6, 9, 12, 15, 18, 21, 24
A		
A	UEMURA, Y. et al., Grammatically Modeling and Predicting RNA Secondary Structures. Genome Informatics Series, Proceedings Genome Informatics Workshop, 11 December, 1995 (11.12.95), No.6, pages 67 to 76	1-24

☐ Further documents are listed in the continuation of Box C. ☐ See patent family annex.

* Special categories of cited documents:	"I" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
"A" document defining the general state of the art which is not considered to be of particular relevance	"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
"E" earlier document but published on or after the international filing date	"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	"&" document member of the same patent family
"O" document referring to an oral disclosure, use, exhibition or other means	
"P" document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search 17 February, 2003 (17.02.03)	Date of mailing of the international search report 04 March, 2003 (04.03.03)
---	---

Name and mailing address of the ISA/
Japanese Patent Office

Authorized officer

Facsimile No.

Telephone No.

A. 発明の属する分野の分類 (国際特許分類 (IPC))		
Int. Cl ⁷ G06F17/30, C12Q1/68, C12N15/00, G01N33/50		
B. 調査を行った分野		
調査を行った最小限資料 (国際特許分類 (IPC))		
Int. Cl ⁷ G06F17/30, C12Q1/68, C12N15/00, G01N33/50		
最小限資料以外の資料で調査を行った分野に含まれるもの		
日本国実用新案公報 1922-1996年 日本国公開実用新案公報 1971-2003年 日本国登録実用新案公報 1994-2003年 日本国実用新案登録公報 1996-2003年		
国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)		
JICSTファイル (JOIS), WPI, INSPEC (DIALOG)		
C. 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求の範囲の番号
X	吉沢 他, RNA二次構造予測支援システムの開発と実装, 人工知能学会人工知能基礎論研究会 (第32回) 資料, 1998. 03. 26, p. 29-35 特に, p. 31	1, 2, 4, 5, 7, 8, 10, 11, 13, 14, 16, 17, 19, 20, 22, 23
A		3, 6, 9, 12, 15, 18, 21, 24
<input checked="" type="checkbox"/> C欄の続きにも文献が列挙されている。 <input type="checkbox"/> パテントファミリーに関する別紙を参照。		
* 引用文献のカテゴリー 「A」 特に関連のある文献ではなく、一般的技術水準を示すもの 「E」 国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの 「L」 優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す) 「O」 口頭による開示、使用、展示等に言及する文献 「P」 国際出願日前で、かつ優先権の主張の基礎となる出願日の後に公表された文献 「T」 国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの 「X」 特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの 「Y」 特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの 「&」 同一パテントファミリー文献		
国際調査を完了した日	17. 02. 03	国際調査報告の発送日
国際調査機関の名称及びあて先 日本国特許庁 (ISA/JP) 郵便番号 100-8915 東京都千代田区霞が関三丁目4番3号		特許庁審査官 (権限のある職員) 高瀬 勤 電話番号 03-3581-1101 内線 3597

C (続き). 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求の範囲の番号
A	UEMURA, Y. et al. Grammatically Modeling and Predicting RNA Secondary Structures. Genome Informatics Series, Proceedings Genome Informatics Workshop, 1995. 12. 11, No. 6, p. 67-76	1-24